

SINGLE-CELL RNA SEQUENCING WITH WATERFALL REVEALS MOLECULAR
CASCADES UNDERLYING ADULT HIPPOCAMPAL NEUROGENESIS

by
Jaehoon Shin, MD

A dissertation submitted to Johns Hopkins University in conformity with the
requirements for the degree of Doctor of Philosophy

Baltimore, Maryland
February, 2016

© Jaehoon Shin 2016

All Rights Reserved

Abstract

Somatic stem cells contribute to tissue ontogenesis, homeostasis, and regeneration through sequential processes. Systematic molecular analyses of stem cells and their development are challenging because classic approaches cannot resolve cellular heterogeneity or capture developmental dynamics. Here we provide a comprehensive resource of single-cell transcriptomes of adult hippocampal quiescent neural stem cells (qNSCs) and their immediate progeny. We further developed Waterfall, a bioinformatic suite, to statistically quantify single-cell gene expression along de novo reconstructed continuous developmental trajectory. Our study reveals molecular signatures of qNSCs, characterized by high-niche signaling and low-protein translation capacities. Our analyses further delineate molecular cascades underlying adult qNSC activation and neurogenesis initiation, exemplified by decreased extrinsic signaling capacity, primed translational machinery, and switches in transcription factors, metabolism, and energy sources. Together, our study reveals the molecular continuum underlying adult neurogenesis and illustrates how Waterfall can be used for single-cell omics analyses of various continuous biological processes.

Advisor: Dr. Hongjun Song, PhD

Reader: Dr. Harry C. Dietz, MD

Acknowledgments

I would like to thank the faculty, trainees, and staff of Song laboratory at Johns Hopkins for their guidance and mentorship, specifically, Dr. Song, who inspired me to ask challenging questions and become an innovator in the field.

Table of Contents

Abstract.....	ii
Acknowledgments	iii
Table of Contents.....	iv
List of Tables	v
List of Figures	vi
Chapter 1. Introduction.....	1
Chapter 2. Single cell RNA sequencing and Waterfall Analysis of Adult Neural Stem Cells	4
Chapter 3. In Vivo Molecular Dynamics of Adult Neural Stem Cells.	18
Chapter 4. Discussion	36
Chapter 5. Single cell analysis with Waterfall	41
Chapter 6. Experimental procedures.....	86
References.....	92
Curriculum Vitae	95

List of Tables

Table 2. 1. Summary of Sequencing statistics.....	9
Table 2. 2. List of literatures consistent with the predicted molecular dynamics from this study.....	10
Table 2. 3. Adult NPC-enriched gene list and validation based on the Allen Brain RNA in situ hybridization database.....	11
Table 3. 1. 1000 UP genes and 1000 DOWN genes and their Spearman correlation coefficient to pseudotime.	24
Table 3. 2. List of UP TFs and DOWN TFs.	26
Table 5. 1. List of datasets to which Waterfall was applied.	56
Table 5. 2 . Comparison between static and dynamic heterogeneity.	57

List of Figures

Figure 2. 1. Single-cell transcriptomes of adult neural stem cells and their immediate progeny.....	12
Figure 2. 2. Single-cell RNA-seq of labeled precursor cells from the adult mouse dentate gyrus....	13
Figure 2. 3. Waterfall for analyzing single-cell data from continuous in vivo process.	14
Figure 2. 4. Reproducibility and orientation of the developmental trajectory on the PCA plot.....	15
Figure 2. 5. Validation for Waterfall predictions for early adult neurogenesis.....	16
Figure 2. 6. Waterfall algorithms.	17
Figure 3. 1. Molecular cascade underlying adult quiescent neural stem cell activation and neurogenesis initiation.....	27
Figure 3. 2. Pseudotime profiles of sample transcription factors.	28
Figure 3. 3. Functional characterization of UP1000 genes and DOWN1000 genes.	29
Figure 3. 4. Cellular compartment analysis for UP and DOWN genes.....	30
Figure 3. 5. Sequential molecular dynamics during adult neural stem cell activation and neurogenesis.....	31
Figure 3. 6. Representative pseudotime profiles of key genes related to niche signaling, cell cycle progression, and energy metabolism.	32
Figure 3. 7. Schematic summary of molecular signatures of quiescent adult neural stem cells and molecular cascades underlying their activation and neurogenesis.....	34
Figure 3. 8. Independent validation for GO entity enrichment test for up-regulated and down- regulated genes during quiescent stem cell activation and neurogenesis.....	35
Figure 5. 1. Waterfall pipeline consists of three key steps of Pre-processing, Reconstruction, and Gene expression analysis.	58
Figure 5. 2. Schematic comparison of the complexity between single cell qPCR or mass cytometry.	59

Figure 5. 3. Hidden Markov model (HMM). We predicted underlying states from gene expression (TPM) over pseudotime progression.	60
Figure 5. 4. Precision of CFPnuc+ single cell isolation (A) and minimal cDNA amplification related 3' bias (B).	61
Figure 5. 5. Marker based cell type screening. Single cells expressing known markers for pericytes, oligodendrocyte (OL), oligodendrocyte progenitor cells (OPC) were excluded for the further analyses.	62
Figure 5. 6. Preprocessing step of Waterfall - Unsupervised clustering (A), principal component analysis (B) and marker expression patterns (C) to help orient of data.	63
Figure 5. 7. Potential trajectories of adult neurogenesis dataset- route definition using minimum spanning tree with k-means clustering (A), the major neurogenic trajectory (B) and alternative routes to consider (C).	64
Figure 5. 8. Trajectory reconstruction step of Waterfall. The trajectory was built using MST connecting the five vertices from k-means clustering (A).	65
Figure 5. 9. Representative gene expression dynamics predicted by Waterfall. Shown are scatterplots of single cell gene expression levels.	66
Figure 5. 10. Monocle-derived trajectory using top 1000 highly expressed genes.	67
Figure 5. 11. The quantitative test for gradual transcriptomic transition of single cell trajectories. ..	68
Figure 5. 12. Expression profiles of key neurogenic genes from Monocle (A) and Waterfall (B).	69
Figure 5. 13. Preprocessing steps of Waterfall: Unsupervised clustering (A), marker gene expression profiles (B), Waterfall trajectory determination without prior temporal information (C); and Monocle trajectory with prior temporal information as a comparison (D).	70
Figure 5. 14. Pseudotime reconstruction of selected route of in vitro HSMM differentiation. MST connected k-means to generate trajectory (A).	71
Figure 5. 15. Gene expression dynamics over the reconstructed HSMM differentiation process.	72
Figure 5. 16. Comparison between Waterfall and Monocle for in vitro myogenesis dataset.	73

Figure 5. 17. Preprocessing step of Waterfall for embryonic lung development: including unbiased clustering (A), PCA plot using selected genes from the four major PC axes (B), and marker gene expression profiles to orient and define a single route of interest (C).....	74
Figure 5. 18. PCA plot after removing Clara/ciliated cell group.	75
Figure 5. 19. Reconstructing continuous trajectory and assigning pseudotime for AT2 (A) and AT1 differentiation process (B).	76
Figure 5. 20. Waterfall gene expression prediction profile for AT2 markers (A) and AT1 markers (B).	77
Figure 5. 21. Cluster based differential analysis adopted from Figure 2A of Treutlein et al. (Treutlein et al., 2014).	78
Figure 5. 22. Waterfall trajectory reconstruction for the “Sample C” of Mass cytometry B cell development dataset. Vertices from the k-means clustering were connected by MST (left).	79
Figure 5. 23. Waterfall predicted gene expression profiles along the B cell developmental trajectory from the “Sample B” (A and B) and the Wanderlust predicted profile (C).	80
Figure 5. 24. Preprocessing of Waterfall analysis for synthetic dataset I.....	81
Figure 5. 25. Preprocessing of Waterfal analysis for the synthetic dataset II.	82
Figure 5. 26. Waterfall reconstructed trajectory and pseudotime assignment for synthetic dataset I.	83
Figure 5. 27. Waterfall reconstructed trajectory and pseudotime assignment for synthetic dataset II.	84
Figure 5. 28. Gene expression profiles for synthetic dataset I, predicted by Waterfall.	85

Chapter 1. Introduction¹

¹ This chapter is based on Shin, J., et al. (2015). "Single-cell RNA-seq with waterfall reveals molecular cascades underlying adult neurogenesis." Cell Stem Cell **17**(3): 360-372.

In discrete regions of the adult mammalian brain, quiescent neural stem cells (qNSCs) continuously generate new neurons through a recurrent process involving quiescent to active state transitions, cell cycle entry and neuronal fate specification (Ming and Song 2011). Understanding molecular mechanisms underlying adult NSC regulation and neurogenesis will not only advance our knowledge of neural development and plasticity, but enable new approaches for regenerative medicine and treatment of brain disorders. Mechanistic analysis of stem cell biology requires comprehensive quantification of molecular properties, such as gene expression. In contrast to traditional approaches targeting individual candidate genes, transcriptome profiling through RNA sequencing (RNA-seq) provides an unbiased and quantitative proxy for molecular features of cellular states. Such a blueprint may reveal unexpected features of NSC biology, generate hypotheses for functional analysis, and lead to novel strategies to manipulate neurogenesis processes.

Classic approaches for molecular characterization of somatic stem cell behavior use population-based readouts at a few time points along development, which faces two major challenges: resolving cellular heterogeneity and capturing developmental dynamics. Adult stem cells constitute a minor population within complex tissues, intermingled with their progeny at different developmental stages and supporting cells. They switch among different states, such as quiescence and activation (Li and Clevers 2010), and thus exhibit significant cellular and molecular differences even upon prospective isolation via fluorescence-activated cell sorting or genetic labeling (Lu, Neff et al. 2011, Codega, Silva-Vargas et al. 2014). Furthermore, snapshots of molecular composition at selected time points are not sufficient to understand the dynamic nature of stem cell development.

Single-cell RNA-seq generates gene expression profiles at the resolution of an individual cell and has thus far revealed molecular profiles of cell types that were not previously recognized at the population level (Stegle, Teichmann et al. 2015). Single-cell RNA-seq has not yet been widely adopted for adult somatic stem cell studies due to technical difficulties in obtaining individual stem cells from complex tissues. Further, the stochastic nature of gene expression in individual cells (Novick and Weiner 1957, Raj, Peskin et al. 2006, Muramoto, Cannon et al. 2012) may lead to overestimation of cellular heterogeneity and requires a new approach for statistical quantification. And for biological systems with only a few known markers, current approaches are

not sufficient to map out a developmental trajectory at high resolution from single-cell datasets (Bendall, Davis et al. 2014, Trapnell, Cacchiarelli et al. 2014).

Despite recent advances in acquiring snapshots of transcriptomes, epigenomes, and proteomes from individual cells, a remaining hurdle is the lack of methodology to identify molecular state transitions over a developmental continuum. Cells are destroyed during acquisition of omic data, so the same cell can't be tracked over time. Here we developed a conceptually different approach, analogous to the "shot-gun" method used in the human genome project characterized by parallel sequencing and bioinformatic reconstruction. We focused on the narrow time window of adult qNSC activation and neurogenesis initiation. Using the *Nestin-CFP^{huc}* transgenic genetic labeling system, we produced single-cell transcriptomes from a mixed population of precursor cells at different developmental stages. We then developed a bioinformatic pipeline named Waterfall to reconstruct continuous biological processes at single-cell resolution using adult neurogenesis as our model, and applied this methodology to other stem cell datasets.

Chapter 2. Single cell RNA sequencing and Waterfall Analysis of Adult Neural Stem Cells²

² This chapter is based on Shin, J., et al. (2015). "Single-cell RNA-seq with waterfall reveals molecular cascades underlying adult neurogenesis." Cell Stem Cell **17**(3): 360-372.

Single-cell RNA-seq of Neural Precursor Cells from the Adult Mouse

Dentate Gyrus

In the adult dentate gyrus, radial glia-like qNSCs give rise to new neurons via a sequential process of activation, proliferation and generation of intermediate precursor cells (IPCs; Figure 2. 1A) (Ming and Song 2011). To elucidate detailed molecular dynamics during initial phases of adult neurogenesis in vivo, we used a transgenic mouse line that expresses nucleus-localized cyan fluorescent protein (CFP) under the Nestin promoter (*Nes-CFP^{nuc}*) (Encinas, Vaahtokari et al. 2006), which labels the majority of NSCs and their immediate progeny (collectively named NPCs; Figure 2. 2A). The SMART-seq protocol (Ramskold, Luo et al. 2012) was modified by adding DNase I treatment step to remove genomic DNA for single-cell cDNA amplification (Figure 2. 2A). In total, we performed single-cell RNA-seq for 142 CFP^{nuc+} and 26 CFP^{nuc-} single cells (Table 2. 1). Total RNA from wild-type adult mouse dentate gyri was serially diluted to 3 pg and processed in parallel for comparison.

We achieved, on average, 87% mapping onto annotated genes (Table 2. 1). Sequencing reads were evenly distributed throughout the whole span of transcripts with 3' bias comparable to recent studies (See Chapter 5. Single cell analysis with Waterfall). Correlation analyses of RNA samples from different batches indicated minimal technical fluctuation during cDNA amplification or across batches compared to significant biological heterogeneity among single-cell transcriptomes (Figure 2. 2B; Table 2. 1). Universally expressed genes, such as β -Actin/Actb, Gapdh, or Ubiquitin B/Ubb, showed even expression patterns across all individual cells, whereas known NSC markers Gfap and Sox2, or early IPC (eIPC) markers Tbr2/Eomes and Sox11, were expressed in subsets of cells (Figure 2. 1C; Table 2. 2).

Nes-CFP^{nuc} also labeled a small percentage of non-NPCs in the adult dentate gyrus (Figure 2. 2A and C). CFP transcript levels were multiple orders of magnitude higher in CFP^{nuc+} cells compared to CFP^{nuc-} cells or diluted dentate RNA (Figure 2. 1B). We excluded cells that exhibited markedly different transcriptomic profiles from the majority of the CFP^{nuc+} population or were clearly identifiable as non-NPCs, such as oligodendrocyte progenitor cells or pericytes (Figure 2. 2C). By differential expression analysis, we identified the top 35 genes enriched in CFP^{nuc+} cells, which included known NSC markers, Blbp, Spot14/Thrsp, Sox9 and GLAST/Slc1a3 (Table 2. 2 and Table 2. 3). In total, 31 out of 35 top genes exhibited SGZ-enriched expression patterns and/or

were known NPC genes ($p = 6.1 \times 10^{-40}$; hypergeometric test; Table 2. 3). These results provided initial validation of our approach.

Waterfall: Analyzing Single-Cell Datasets from Continuous in vivo Processes

We next examined the whole-transcriptome dataset of individual CFP^{nuc+} NPCs. Unsupervised hierarchical clustering analysis resulted in two super-groups with six sub-groups (Figure 2. 3A). Notably, these six CFP^{nuc+} groups were not clearly segregated on the principal component analysis (PCA) plot (Figure 2. 3B), which was consistent over different batches of sequencing runs with multiple biological replicates (Figure 2. 4A; Table 2. 1). The continuous trajectory called for a new approach not relying on segmentation into a few groups of cell clusters. We could not use currently available single-cell analysis software, such as Monocle (Trapnell, Williams et al. 2010) or Wanderlust (Bendall, Davis et al. 2014), for our system due to the lack of sufficient prior information, such as temporal delineators or a robust set of specific markers (See Chapter 5. Single cell analysis with Waterfall). We thus developed a more generally applicable pipeline of algorithms to perform unbiased statistical analyses of multidimensional single-cell datasets from continuous biological processes. We collectively named the suite of algorithms “Waterfall”, which involves three steps: pre-processing, pseudotime reconstruction, and gene expression analysis (Figure 2. 3C; See Chapter 5. Single cell analysis with Waterfall).

Pre-processing defined the trajectory of interest following dimensionality reduction of the data. Unsupervised learning identified six clusters of cells (Figure 2. 3A), which were then labeled S1-5 and SA based on their relative location in a PCA plot (Figure 2. 3B). SA was recognized as a branch by the minimum spanning tree (MST) algorithm (See Chapter 5. Single cell analysis with Waterfall). Although characterizing SA would be interesting (See Chapter 5. Single cell analysis with Waterfall), we focused on the major neurogenic pathway in the current study. The expression profiles of a few known developmental genes were used to orient the most probable trajectory of interest (Figure 2. 3B and Figure 2. 4B).

To reconstruct the chronology, we first determined the most probable route of transcriptomic progression. We performed k-means clustering of single-cell transcriptomes on the PCA plot after excluding SA, followed by constructing a MST trajectory to connect cluster centers (Figure 2. 6A). We then introduced “pseudotime” (Trapnell, Cacchiarelli et al. 2014) to define the relative location

of each cell on the MST trajectory (Figure 2. 6A). Taking the Euclidian distance defined by the whole transcriptomic difference from each cell to the next in pseudotime, we found that the total path length reconstructed by Waterfall was significantly shorter than would result from random ordering of cells (See Chapter 5. Single cell analysis with Waterfall). The pseudotime algorithm reconstructs molecular state transitions of a continuous process by quantifying the gradual divergence of single-cell transcriptomes individually, rather than as members of pre-classified groups.

For gene expression analysis, we developed an algorithm to determine the binary on/high or off/low expression state of each gene along pseudotime in an unbiased fashion using a hidden Markov model (HMM; Figure 2. 6B). Gene expression at the single-cell level is highly stochastic and binary (Novick and Weiner 1957, Raj, Peskin et al. 2006, Muramoto, Cannon et al. 2012). When analyzing continuous processes represented by single-cell transcriptomes in the absence of discrete groups, conventional statistical methods, such as arithmetic mean or t-test, are not appropriate. We adopted HMM to statistically convert stochastic expression patterns of individual genes into binary on/high or off/low states (Figure 2. 6B). The binary gene expression states were then shown as heat maps, which quantified the molecular cascade over time (Figure 2. 3C-D). To discover novel developmentally regulated genes, we correlated gene expression levels with pseudotime and subjected identified genes to gene ontology (GO) analyses (Figure 2. 3C).

Validation for the Reconstructed Adult Neurogenesis Process

We validated molecular dynamics revealed by Waterfall at multiple levels. First, known NSC markers *Gfap* and *Apoe*, and eIPC markers *Sox11* and *Tbr2*, showed non-overlapping expression states over developmental pseudotime (Figure 2. 3D; Table 2. 2). Second, we evaluated in vivo expression of *Aldoc* and *Stmn1*, which have not been previously studied in adult hippocampal neurogenesis. Over pseudotime, *Aldoc* was initially highly expressed (on/high state), then downregulated (off/low state), whereas *Stmn1* was initially off, then upregulated (Figure 2. 5A and C). In the hippocampal dentate subgranular zone (SGZ) of adult *Nes-GFP^{cyto}* mice (Encinas, Michurina et al. 2011), *Aldoc*⁺*GFP*⁺ precursors were almost exclusively *PCNA*⁻ qNSCs, with very few *PCNA*⁺ active NSCs (aNSCs) or IPCs (Figure 2. 5B). In contrast, *Stmn1*⁺*GFP*⁺ precursors were mostly eIPCs and *PCNA*⁺ aNSCs, but not *PCNA*⁻ qNSCs (Figure 2. 5D). Thus, Waterfall accurately

predicted the in vivo expression dynamics of both known and unknown genes. Third, for functional validation, we explored the possibility of genetic labeling of a specific developmental stage based on Waterfall results. Hopx (Hop homeobox) was highly expressed in qNSCs, but downregulated around the transition point from qNPC to aNSC (Figure 2. 5E). Upon a single low-dose tamoxifen injection into a *Hopx-CreER^{T2}* mouse line (Takeda, Jain et al. 2011) crossed with a *mT/mG^{f/f}* reporter line for clonal lineage-tracing (Bonaguidi, Wheeler et al. 2011), almost all labeled precursors at three days post injection were nestin⁺GFAP⁺ qNSCs in the adult SGZ (Figure 2. 5F). By 7 days, we were able to observe GFP-labeled clones that contained both NSCs and their progeny (Figure 2. 5G), indicating self-renewal and differentiation, two hallmarks of NSCs.

Table 2. 1. Summary of Sequencing statistics.

Batch names	Hiseq5	Hiseq6	Hiseq7/8*	Hiseq9	Hiseq10	Total
Run type	100 cycle single end	100 cycle single end	100 cycle single end	50 cycle single end	100 cycle single end	
Total read number	106,765,153	102,075,273	219,505,841	113,311,365	103,066,717	644,724,349
Number samples	19	24	53	34	42	172
Number CFP+ cells	19	19	39	30	35	142
Number CFP- cells	0	3	12	4	7	26
Number bulk RNA samples	0	2	2	0	0	4

Table 2. 2. List of literatures consistent with the predicted molecular dynamics from this study.

Developmental genes	
Spot14/Thrsp	Knobloch, M., Braun, S.M., Zurkirchen, L., von Schoultz, C., Zamboni, N., Arauzo-Bravo, M.J., Kovacs, W.J., Karalay, O., Suter, U., Machado, R.A., et al. (2013). Metabolic control of adult neural stem cell activity by Fasn-dependent lipogenesis. <i>Nature</i> 493, 226-230.
Tbr2/Eomes	Hodge, R.D., Kowalczyk, T.D., Wolf, S.A., Encinas, J.M., Rippey, C., Enikolopov, G., Kempermann, G., and Hevner, R.F. (2008). Intermediate progenitors in adult hippocampal neurogenesis: Tbr2 expression and coordinate regulation of neuronal output. <i>J Neurosci</i> 28, 3707-3717.
Sox11	Haslinger, A., Schwarz, T.J., Covic, M., and Lie, D.C. (2009). Expression of Sox11 in adult neurogenic niches suggests a stage-specific role in adult neurogenesis. <i>The European journal of neuroscience</i> 29, 2103-2114.
Sox9	Scott, C.E., Wynn, S.L., Sesay, A., Cruz, C., Cheung, M., Gomez Gaviro, M.V., Booth, S., Gao, B., Cheah, K.S., Lovell-Badge, R., et al. (2010). SOX9 induces and maintains neural stem cells. <i>Nature neuroscience</i> 13, 1181-1189.
GLAST/Slc1a3	Doetsch, F. (2003). The glial identity of neural stem cells. <i>Nature neuroscience</i> 6, 1127-1134.
ApoE	Li, G., Bien-Ly, N., Andrews-Zwilling, Y., Xu, Q., Bernardo, A., Ring, K., Halabisky, B., Deng, C., Mahley, R.W., and Huang, Y. (2009). GABAergic interneuron dysfunction impairs hippocampal neurogenesis in adult apolipoprotein E4 knockin mice. <i>Cell Stem Cell</i> 5, 634-645.
Id3	Niola, F., Zhao, X., Singh, D., Castano, A., Sullivan, R., Lauria, M., Nam, H.S., Zhuang, Y., Benezra, R., Di Bernardo, D., et al. (2012). Id proteins synchronize stemness and anchorage to the niche of neural stem cells. <i>Nature cell biology</i> 14, 477-487.
Nr2e1/Tlx	Zhang, C.L., Zou, Y., He, W., Gage, F.H., and Evans, R.M. (2008). A role for adult TLX-positive neural stem cells in learning and behaviour. <i>Nature</i> 451, 1004-1007.
Hes1	Hatakeyama, J., Bessho, Y., Katoh, K., Ookawara, S., Fujio, M., Guillemot, F., and Kageyama, R. (2004). Hes genes regulate size, shape and histogenesis of the nervous system by control of the timing of neural stem cell differentiation. <i>Development</i> 131, 5539-5550.
SoxC (Sox4 and Sox11)	Mu, L., Berti, L., Masserdotti, G., Covic, M., Michaelidis, T.M., Doberauer, K., Merz, K., Rehfeld, F., Haslinger, A., Wegner, M., et al. (2012). SoxC transcription factors are required for neuronal differentiation in adult hippocampal neurogenesis. <i>J Neurosci</i> 32, 3067-3080.
Foxg1	Shen, L., Nam, H.S., Song, P., Moore, H., and Anderson, S.A. (2006). FoxG1 haploinsufficiency results in impaired neurogenesis in the postnatal hippocampus and contextual memory deficits. <i>Hippocampus</i> 16, 875-890.
Insm1	Duggan, A., Madathany, T., de Castro, S.C., Gerrelli, D., Guddati, K., and Garcia-Anoveros, J. (2008). Transient expression of the conserved zinc finger gene INSM1 in progenitors and nascent neurons throughout embryonic and adult neurogenesis. <i>The Journal of comparative neurology</i> 507, 1497-1520.
Insm1	Rosenbaum, J.N., Duggan, A., and Garcia-Anoveros, J. (2011). Insm1 promotes the transition of olfactory progenitors from apical and proliferative to basal, terminally dividing and neuronogenic. <i>Neural development</i> 6, 6.
Tcf12	Uittenbogaard, M., and Chiaramello, A. (2002). Expression of the bHLH transcription factor Tcf12 (ME1) gene is linked to the expansion of precursor cell populations during neurogenesis. <i>Brain research Gene expression patterns</i> 1, 115-121.
Nf1b	Ninkovic, J., Steiner-Mezzadri, A., Jawerka, M., Akinci, U., Masserdotti, G., Petricca, S., Fischer, J., von Holst, A., Beckers, J., Lie, C.D., et al. (2013). The BAF complex interacts with Pax6 in adult neural progenitors to establish a neurogenic cross-regulatory transcriptional network. <i>Cell Stem Cell</i> 13, 403-418.
Dbx2	Pierani, A., Brenner-Morton, S., Chiang, C., and Jessell, T.M. (1999). A sonic hedgehog-independent, retinoid-activated pathway of neurogenesis in the ventral spinal cord. <i>Cell</i> 97, 903-915.
Nr3c1	Sundberg, M., Savola, S., Hienola, A., Korhonen, L., and Lindholm, D. (2006). Glucocorticoid hormones decrease proliferation of embryonic neural stem cells through ubiquitin-mediated degradation of cyclin D1. <i>J Neurosci</i> 26, 5402-5410.
Id4	Bedford, L., Walker, R., Kondo, T., van Cruchten, I., King, E.R., and Sablitzky, F. (2005). Id4 is required for the correct timing of neural differentiation. <i>Developmental biology</i> 280, 386-395.
Hmgb1	Merenmies, J., Pihlaskari, R., Laitinen, J., Wartiovaara, J., and Rauvala, H. (1991). 30-kDa heparin-binding protein of brain (amphoterin) involved in neurite outgrowth. Amino acid sequence and localization in the filopodia of the advancing plasma membrane. <i>The Journal of biological chemistry</i> 266, 16722-16729.
N-myc	Knoepfler, P.S., Cheng, P.F., and Eisenman, R.N. (2002). N-myc is essential during neurogenesis for the rapid expansion of progenitor cell populations and the inhibition of neuronal differentiation. <i>Genes & development</i> 16, 2699-2712.
Mxd3	Shimada, Y., Kuroyanagi, J., Zhang, B., Ariyoshi, M., Umemoto, N., Nishimura, Y., and Tanaka, T. (2014). Downregulation of Max dimerization protein 3 is involved in decreased visceral adipose tissue by inhibiting adipocyte differentiation in zebrafish and mice. <i>International journal of obesity</i> 38, 1053-1060.
Zeb2	Denecker, G., Vandamme, N., Akay, O., Koludrovic, D., Taminau, J., Lemeire, K., Gheldof, A., De Craene, B., Van Gele, M., Brochez, L., et al. (2014). Identification of a ZEB2-MITF-ZEB1 transcriptional network that controls melanogenesis and melanoma progression. <i>Cell death and differentiation</i> 21, 1250-1261.
ZT3/Zip	Polimeni, M., Giorgi, S., De Gregorio, L., Dragani, T.A., Molinaro, M., Cossu, G., and Bouche, M. (1996). Differentiation dependent expression in muscle cells of ZT3, a novel zinc finger factor differentially expressed in embryonic and adult tissues. <i>Mechanisms of development</i> 54, 107-117.
Tsc22d3/Gilz	Bruscoli, S., Donato, V., Velardi, E., Di Sante, M., Migliorati, G., Donato, R., and Riccardi, C. (2010). Glucocorticoid-induced leucine zipper (Gilz) and long Gilz inhibit myogenic differentiation and mediate anti-myogenic effects of glucocorticoids. <i>The Journal of biological chemistry</i> 285, 10385-10396.
Hopx	Takeda, N., Jain, R., LeBoeuf, M.R., Wang, Q., Lu, M.M., and Epstein, J.A. (2011). Interconversion between intestinal stem cell populations in distinct niches. <i>Science</i>
Hopx	Takeda, N., Jain, R., LeBoeuf, M.R., Padmanabhan, A., Wang, Q., Li, L., Lu, M.M., Millar, S.E., and Epstein, J.A. (2013). Hopx expression defines a subset of multipotent hair follicle stem cells and a progenitor population primed to give rise to K6+ niche cells. <i>Development</i> 140, 1655-1664.
Developmental signaling pathways	
Insulin/IGF2 signaling	Bracko, O., Singer, T., Aigner, S., Knobloch, M., Winner, B., Ray, J., Clemenson, G.D., Jr., Suh, H., Couillard-Despres, S., Aigner, L., et al. (2012). Gene expression profiling of neural stem cells and their neuronal progeny reveals IGF2 as a regulator of adult hippocampal neurogenesis. <i>J Neurosci</i> 32, 3376-3387.
NMDA signaling	Cameron, H.A., McEwen, B.S., and Gould, E. (1995). Regulation of adult neurogenesis by excitatory input and NMDA receptor activation in the dentate gyrus. <i>J Neurosci</i> 15, 4687-4692.
Wnt signaling	Lie, D.C., Colamarino, S.A., Song, H.J., Desire, L., Mira, H., Consiglio, A., Lein, E.S., Jessberger, S., Lansford, H., Dearie, A.R., et al. (2005). Wnt signalling regulates adult hippocampal neurogenesis. <i>Nature</i> 437, 1370-1375.
BMP signaling	Lim, D.A., Tramontin, A.D., Trevejo, J.M., Herrera, D.G., Garcia-Verdugo, J.M., and Alvarez-Buylla, A. (2000). Noggin antagonizes BMP signaling to create a niche for adult neurogenesis. <i>Neuron</i> 28, 713-726.
Bdnf signaling	Pencea, V., Bingaman, K.D., Wiegand, S.J., and Luskin, M.B. (2001). Infusion of brain-derived neurotrophic factor into the lateral ventricle of the adult rat leads to new neurons in the parenchyma of the striatum, septum, thalamus, and hypothalamus. <i>J Neurosci</i> 21, 6706-6717.
GABA signaling	Song, J., Zhong, C., Bonaguidi, M.A., Sun, G.J., Hsu, D., Gu, Y., Meletis, K., Huang, Z.J., Ge, S., Enikolopov, G., et al. (2012). Neuronal circuitry mechanism regulating adult quiescent neural stem-cell fate decision. <i>Nature</i> 489, 150-154.
FGF signaling	Yoshimura, S., Takagi, Y., Harada, J., Teramoto, T., Thomas, S.S., Waerber, C., Bakowska, J.C., Breakefield, X.O., and Moskowitz, M.A. (2001). FGF-2 regulation of neurogenesis in adult hippocampus after brain injury. <i>Proceedings of the National Academy of Sciences of the United States of America</i> 98, 5874-5879.
Notch signaling	Pleasure, S.J., Collins, A.E., and Lowenstein, D.H. (2000). Unique expression patterns of cell fate molecules delineate sequential stages of dentate gyrus development. <i>J Neurosci</i> 20, 6095-6105.

Table 2. 3. Adult NPC-enriched gene list and validation based on the Allen Brain RNA in situ hybridization database.

Uncorrected p values (Wilcoxon rank sum)	Corrected p values (Benjamini & Hochberg)	log2FD	Covar1	Covar2	TPM_CFP+, NPC	TPM_CFP-, neuronal	Gene_name	Allen atlas SGZ
3.12E-06	1.35E-03	5.513398304	0.866476274	2.05026401	593.1929231	12.98666667	Ptprz1	yes
8.45E-06	1.65E-03	7.148029451	1.032782959	1.697798523	6861.155154	48.37555556	Fabp7	yes
8.45E-06	1.65E-03	3.049352557	0.64142578	0.782681052	10472.36431	1265.022222	Dbl	yes
1.21E-05	1.73E-03	4.503347823	0.672894701	1.245446016	4355.766846	192.0533333	Slc1a3	yes
1.33E-05	1.73E-03	4.101642441	0.858969222	2.103493065	2185.322231	127.2911111	Cd9	yes
2.42E-05	2.90E-03	8.591404217	1.088297417	0.580234811	2893.316692	7.501111111	Thrsp	yes
4.11E-05	3.95E-03	4.349536335	0.802908031	1.792381448	1322.353615	64.86444444	Ppap2b	yes
9.08E-05	0.006641362	2.396843788	0.831080278	0.812084385	991.1899231	188.2066667	Hnrnpa2b1	yes
1.03E-04	6.76E-03	4.981848298	1.017756909	2.200124636	1228.779846	38.88555556	2610017109Rik	yes
1.04E-04	0.006755994	3.205671563	0.862564264	1.840240225	1609.048538	174.4077778	Mmd2	ambiguous
1.09E-04	6.79E-03	2.795301582	0.82375671	1.506968619	658.4336923	94.85111111	Sox9	yes
2.06E-04	0.011947516	8.087390404	1.215233001	1.412191616	851.6194615	3.131111111	B2m	ambiguous
2.07E-04	0.011947516	2.150388132	0.737266688	1.170528541	2389.032846	538.1344444	Ndufc2	yes
2.31E-04	0.012839875	2.36910321	0.708327211	0.584880353	22575.10131	4369.77	Cst3	yes
2.56E-04	0.013731082	3.143811704	1.085206643	2.045264833	429.9273077	48.64222222	Hopx	yes
2.98E-04	0.013950151	3.996824191	0.912255588	1.516816198	4339.574538	271.8211111	Gstm1	ambiguous
2.98E-04	0.013950151	2.89511646	1.168269207	2.910818862	924.8958462	124.33	Cnn3	yes
3.11E-04	0.013950151	3.718420308	1.134249069	1.622670181	317.0097692	24.08333333	Dcll1	no
3.13E-04	0.013950151	9.318514394	1.669742053	1.559022993	439.1374615	0.687777778	Cdk4	yes
3.84E-04	0.016151265	9.906787182	2.415999817	2.471118293	117.3249231	0.122222222	2810055G20Rik	yes
4.01E-04	0.016430184	4.777984082	1.17486707	1.444870629	161.505	5.886666667	Pbx1	no
4.32E-04	0.016846943	5.77146291	1.811156967	2.116127144	109.9156154	2.012222222	Pitpnc1	no
4.33E-04	0.016846943	8.277186919	1.488233106	1.424655853	184.0688462	0.593333333	Eps15	yes
4.54E-04	0.016972126	2.772902073	0.802346253	1.126822052	836.5812308	122.4	Slc1a2	yes
4.58E-04	0.016972126	7.271714896	1.68785924	1.694498596	164.6570769	1.065555556	Pdpn	no
4.73E-04	0.017139429	4.043009581	1.058761678	2.774219046	1191.988846	72.31111111	Mfge8	yes
5.14E-04	0.017783341	2.309309182	0.868634517	1.133779149	1626.681615	328.1944444	Gpm6b	no
5.31E-04	0.017969691	7.76522177	1.4138105	1.804086119	367.6642308	1.69	Naaa	yes
6.15E-04	0.019550025	4.031442716	0.919928845	1.572202299	226.6098462	13.85777778	Gstm3	ambiguous
7.35E-04	0.022902946	7.315668711	1.124599768	1.012553256	1522.623923	9.557777778	Cmtm5	yes
7.66E-04	0.022912004	5.410329823	1.144439386	2.605218371	280.5408462	6.596666667	Ddah1	yes
7.79E-04	0.022912004	6.77820491	1.617221312	2.111178271	386.1103077	3.517777778	Id4	yes
8.07E-04	0.023273379	9.373295493	2.10550752	0.897732718	136.324	0.205555556	Tnc	yes
9.76E-04	0.027150686	1.762032045	0.705377953	0.763674623	4960.751	1462.591111	Ckb	yes

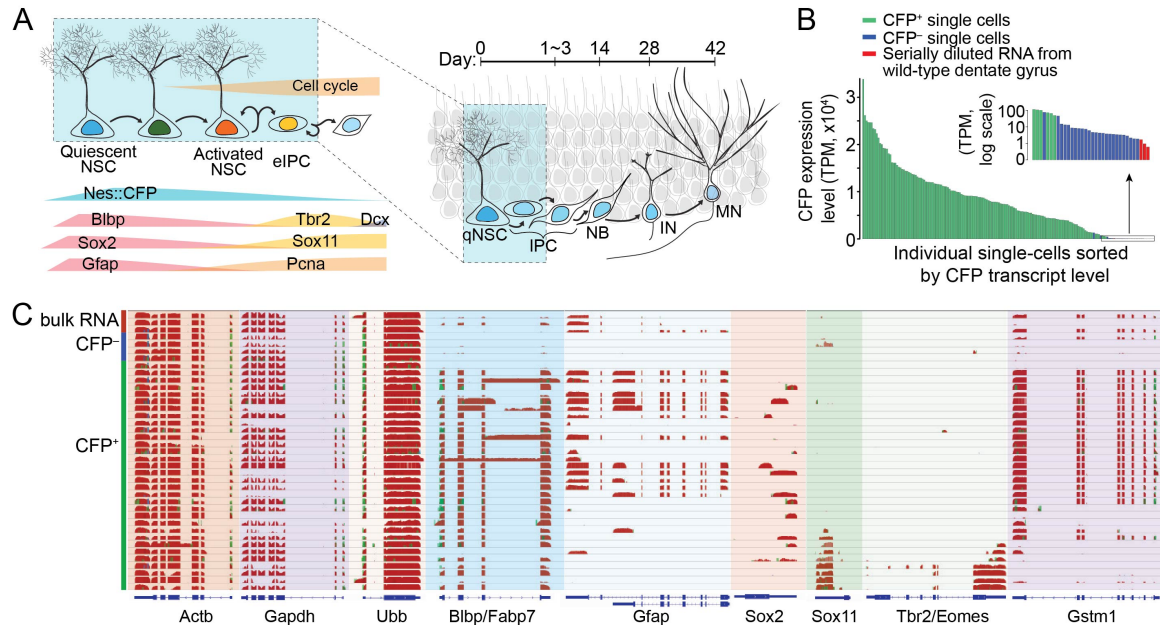


Figure 2. 1. Single-cell transcriptomes of adult neural stem cells and their immediate progeny.

(A) A schematic diagram of the process of adult neurogenesis in the dentate gyrus of the mouse hippocampus. Once quiescent neural stem cells (qNSCs) become activated (aNSCs), they enter cell cycle and generate early intermediate progenitor cells (eIPCs), which in turn give rise to neuroblasts (NB), immature neurons (IN) and then mature neurons (MN). Area highlighted with blue background indicates cell types fluorescently labeled in adult Nestin-CFPnuc animals.

(B) Expression levels of transcript encoding CFP in each single cell and diluted whole dentate RNA samples (TPM, transcripts per million). Inset: enlarged view of CFP transcript levels in logarithmic scale of samples with low abundance of CFP transcript.

(C) Representative coverage profile of diluted total RNA from the whole dentate gyrus, CFP- individual cells, and CFP+ individual cells at selected genomic loci, including house-keeping genes (β -actin/Actb, Gapdh, ubiquitin B/Ubb), known NSC markers (Blbp/Fabp7, Gfap, Sox2), known IPC markers (Sox11, Tbr2/Eomes), and potential new NPC markers (Gstm1).

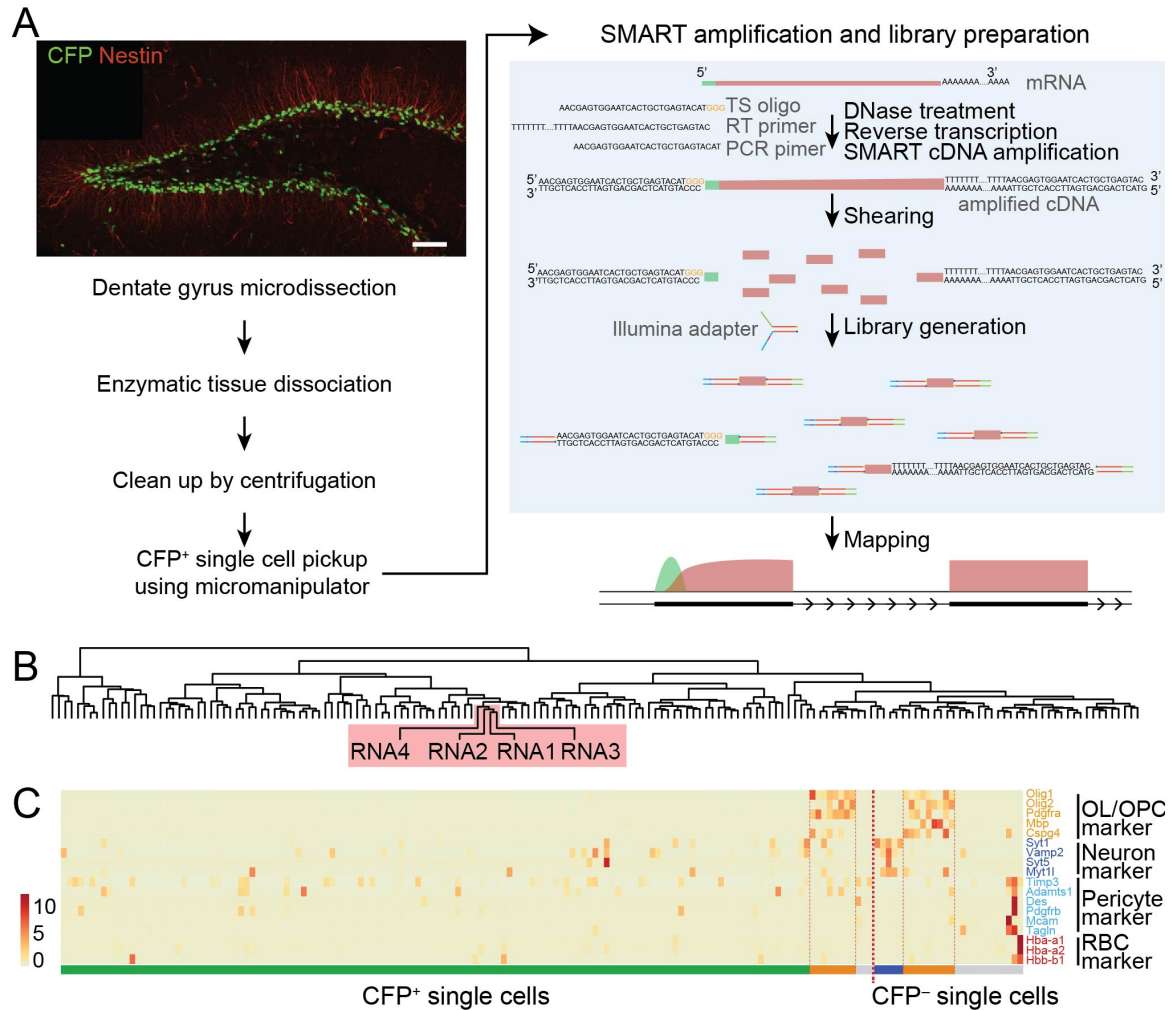


Figure 2. 2. Single-cell RNA-seq of labeled precursor cells from the adult mouse dentate gyrus.

(A) A schematic diagram of experimental procedures for achieving single-cell transcriptomes from fluorescently labeled individual cells isolated from the adult *Nestin-CFP^{nuc}* mouse dentate gyrus. Shown at the top is a sample confocal image of *Nestin-CFP^{nuc}* adult mouse dentate gyrus for CFP and Nestin immunostaining. Scale bar: 100 μ m. Isolation of individual NPCs involved the following steps: (1) microdissection of tissues of interest; (2) dissociation of tissue using papain and DNase I; (3) eliminating cell debris via multiple rounds of mild centrifugation; and (4) picking up fluorescently labeled individual cells using a micromanipulator with a pulled glass micropipette and breaking the tip of the glass pipette into PCR strips. The SMART amplification protocol (Ramskold, Luo et al. 2012) was followed with minor modifications, including DNase I digestion before adding polyT primer. Library generation followed conventional Illumina library preparation after fragmentation of amplified cDNA. TS oligo: Template Switch oligo; RT primer: primer for reverse transcription; SMART: Switching Mechanism At 5' end of RNA Template.

(B) Unsupervised clustering for individual CFP⁺ and CFP⁻ cells, and single-cell equivalent amount (3 pg) of whole-dentate gyrus RNA samples (n = 4). Whole dentate RNA samples are highlighted to show the extensive biological variability of single-cell transcriptomes, compared to negligible technical variability.

(C) Expression heat map of known markers for oligodendrocyte precursor cells (OL)/oligodendrocytes (OL), pericytes, and blood cells. Horizontal bar under the heat map represents putative cell types determined by marker expression profiles. Green: cells with potential NPC lineages; Orange: cells with OPC/OL lineages; Purple: cells with differentiated neuronal characteristics; Gray: uncharacterized outliers or cells with non-neuronal lineages.

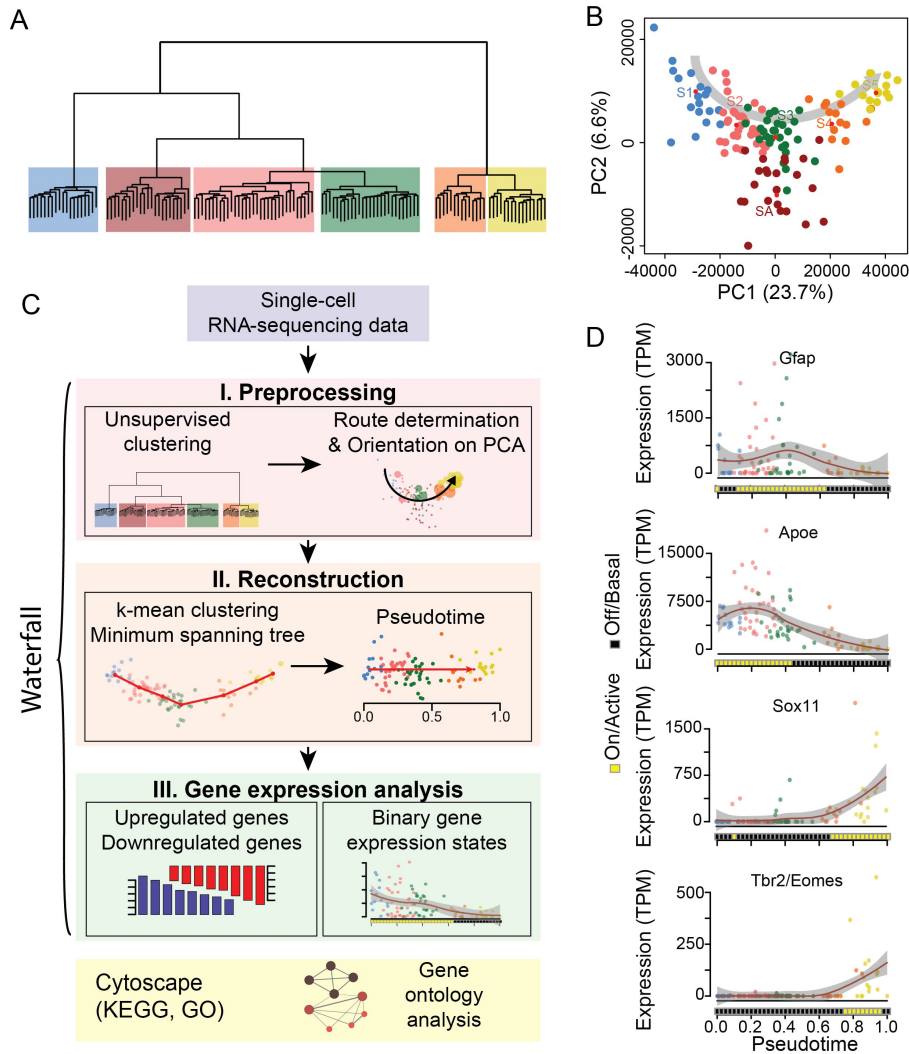


Figure 2. 3. Waterfall for analyzing single-cell data from continuous in vivo process.

(A) Unsupervised clustering analysis of CFP^{nuc+} NPCs resulting in two super-groups with six subgroups. Different groups are color coded in the same fashion in this figure and across all other figures.

(B) Principal component analysis (PCA) plot shows one of the possible linear trajectories of different groups with the exception of SA.

(C) A schematic diagram of multiple components and workflow of Waterfall. Waterfall is a full range of algorithms for processing multi-dimensional single-cell datasets derived from continuous biological processes. Please See Chapter 5. Single cell analysis with Waterfall for more information and Waterfall analyses of other biological systems.

(D) Representative expression profiles of marker genes of adult neurogenesis. Each data point represents the gene expression level of a single cell with color scheme following Figure 2. 3. Data points are fitted with local polynomial regression fitting (red lines) with 95% confidence interval (gray area). HMM-predicted underlying states are represented as black and yellow squares on the bottom of the graphs.

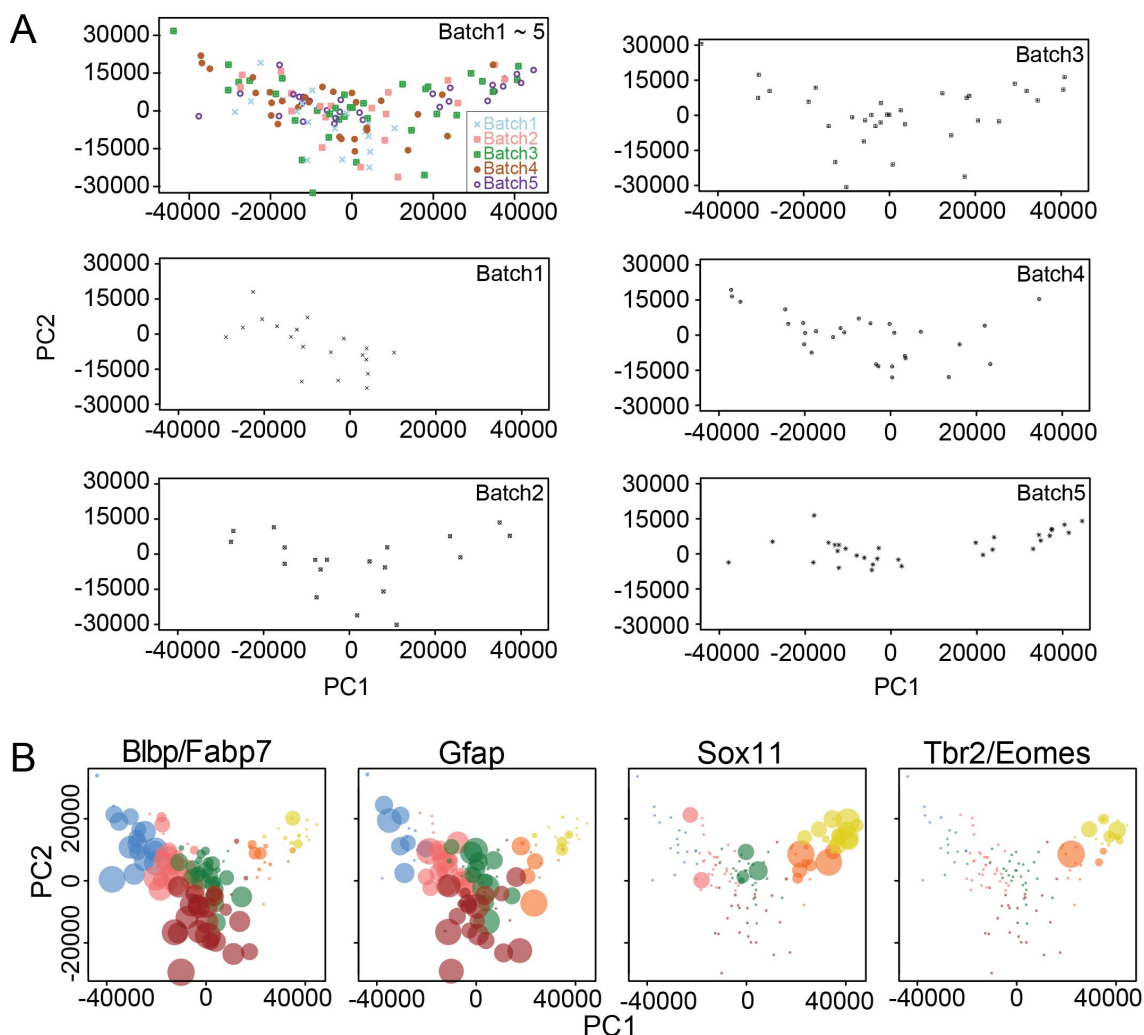


Figure 2. 4. Reproducibility and orientation of the developmental trajectory on the PCA plot

(A) Shown are PCA plots with all Nestin-CFP^{Nuc+} NPCs (top left panel) and with Nestin-CFP^{Nuc+} NPCs in five individual sequencing runs.

(B) Normalized expression levels of known marker genes on the PCA plot, represented by the size of data points. Colors of data points follow the color codes in Figure 2. 3. Notably, known NSC markers *Blbp/Fabp7* and *Gfap* were highly expressed on the left side, whereas eIPC markers *Tbr2/Eomes* and *Sox11* were highly expressed on the right side.

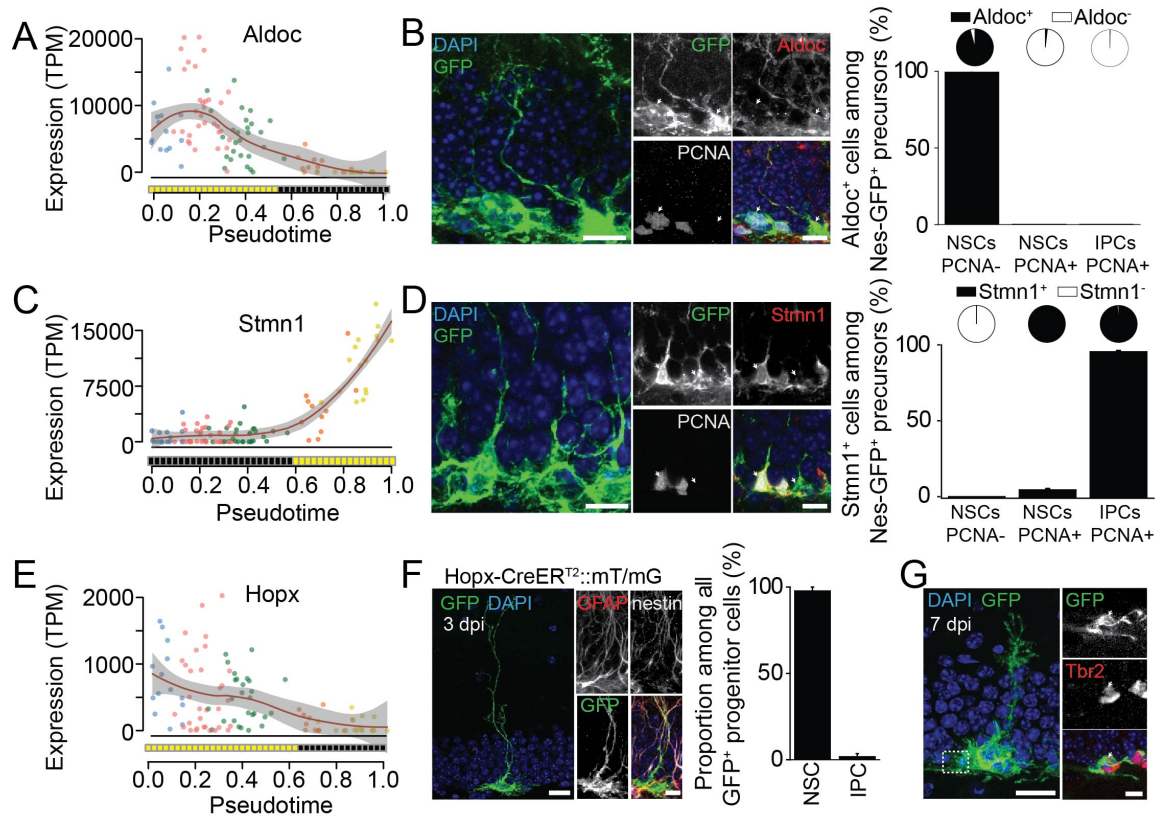


Figure 2. 5. Validation for Waterfall predictions for early adult neurogenesis.

(A-D) Validation of gene expression patterns and on/off binary states of Aldolase C (Aldoc, A) and Stmn1 (C) over pseudotime by immunohistology. Also shown are sample confocal images of GFP, cell proliferation marker PCNA, Aldoc or Stmn1 in the dentate gyrus of adult *Nestin-GFP^{cyco}* mice (left panels) and quantifications (right panels). Values represent mean \pm SEM (n = 3 animals). The pie chart represents the proportion of Aldoc⁺ or Stmn1⁺ cells among each category of the GFP⁺ progenitor population. Scale bars: 20 μ m (left) and 10 μ m (right).

(E-G) Validation of gene expression patterns and on/off binary states of Hopx over pseudotime (E) by genetic labeling and lineage-tracing. Adult *Hopx-CreER^{T2}::mT/mG* mice were injected with a single dose of tamoxifen and examined 3 (F) or 7 days (G) later. Shown in (F) are sample confocal images of GFP, GFAP, Nestin and DAPI. Also shown is quantification of percentages of GFP⁺ cells as NSCs or IPCs. Values represent mean \pm SEM (n = 5 dentate gyri). Shown in (G) is an example of a labeled clone containing an NSC and multiple Tbr2⁺ neuronal progeny. Scale bars: 20 μ m (left) and 10 μ m (right).

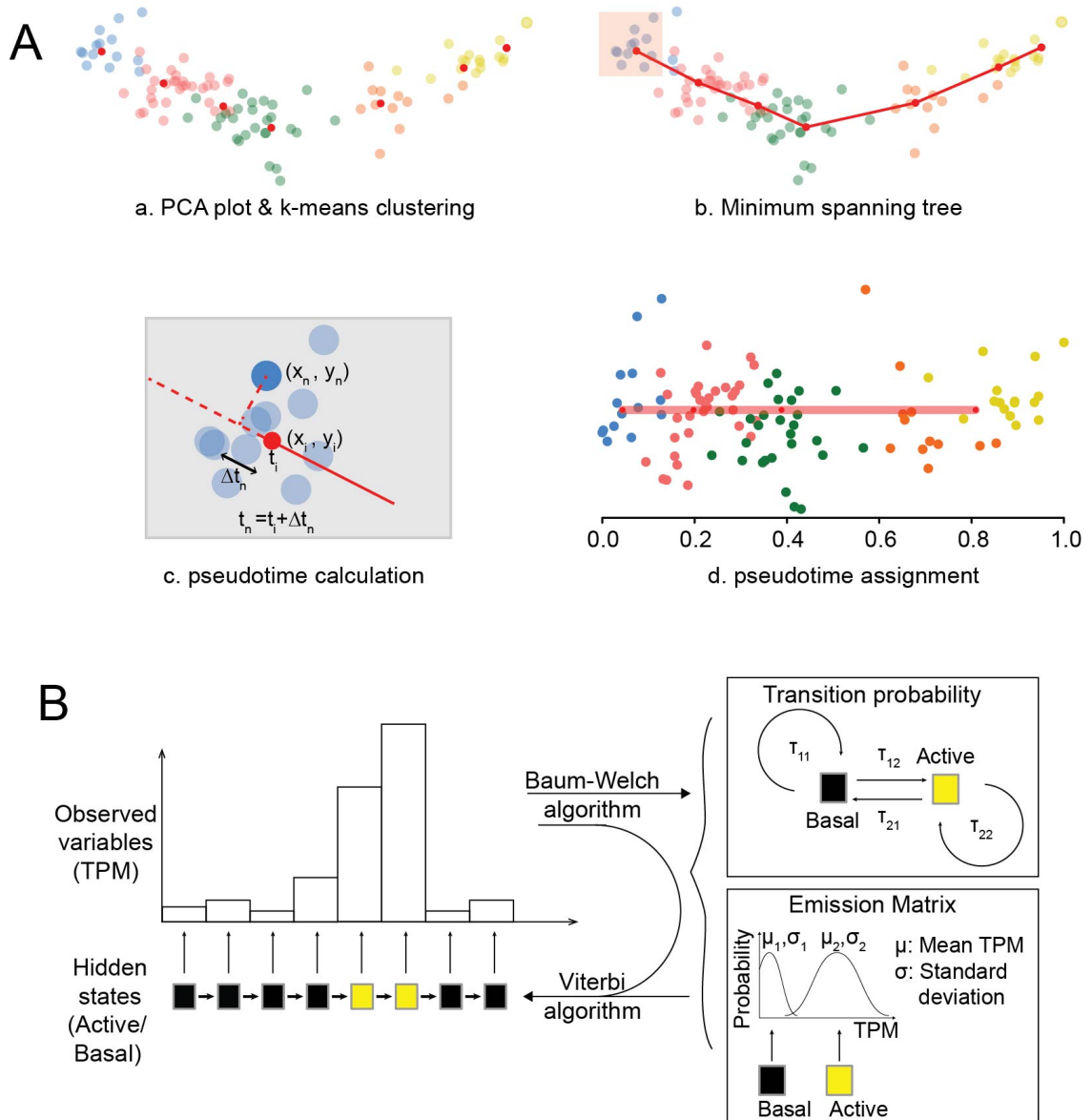


Figure 2. 6. Waterfall algorithms.

(A) An illustration of steps to generate an averaged trajectory and assign pseudotime to each individual cell: (a) Representation of individual cells on the PCA plot with PC1 and PC2 and performing k-means from the PCA data (small red dots); (b) Build a trajectory connecting k-means as minimum spanning trees (MST); (c) Determine the relative location of each cell using orthogonal line to connect each cell to the closest trajectory line; (d) Assign pseudotime for each cell based on its relative location on the trajectory.

(B) An illustration of the approach to predict underlying states from gene expression (TPM) over pseudotime progression. The Baum-Welch algorithm predicts the most likely transition probability and emission matrix from observed variables (TPM). The Viterbi algorithm uses observed variables (TPM) along with the output from the Baum-Welch algorithm to predict hidden On/High and Off/Low gene expression states.

Chapter 3. In Vivo Molecular Dynamics of Adult Neural Stem Cells.³

³ This chapter is based on Shin, J., et al. (2015). "Single-cell RNA-seq with waterfall reveals molecular cascades underlying adult neurogenesis." Cell Stem Cell **17**(3): 360-372.

Transcription Factor Expression during Adult NSC Activation and Neurogenesis

Waterfall allows an unbiased prediction of the relative chronological position of each individual cell and distribution of binary gene expression over the developmental trajectory. To delineate molecular cascades underlying adult qNSC activation and neurogenesis, we generated a list of the top 1000 negatively correlated genes with pseudotime (DOWN¹⁰⁰⁰ genes; Spearman correlation coefficient < -0.13 ; Table 3. 1A), which represent qNSC-enriched genes down-regulated during activation and neurogenesis, as well as the top 1000 positively correlated genes with pseudotime (UP¹⁰⁰⁰ genes; Spearman correlation coefficient > 0.20 ; Table 3. 1B), which represent newly activated genes during qNSC activation and early neurogenesis.

Out of these 2000 genes, we initially focused on transcription factors (TFs). Most known TFs involved in adult neurogenesis were discovered by extrapolating findings from embryonic studies. In contrast, our database provides unbiased genome-wide profiles of TF expression. Systematic analyses of our dataset revealed a total of 41 down-regulated TFs and 42 up-regulated TFs during adult hippocampal neurogenesis (Figure 3. 1 and Figure 3. 2; Table 3. 2).

First, the set of dynamic TFs we identified included known regulators of adult NSCs and neurogenesis, which provided additional validation of our approach. Among DOWN TFs, Sox2, Sox9, Id3, nuclear receptor Nr2e1/Tlx, and Hes1 have been shown to regulate adult NSC maintenance and function. Among UP TFs, SoxC (Sox4 and Sox11), Foxg1, Tbr2, Insm1, Tcf12 and Nfib are critical in proliferative adult NPCs (Table 2. 2).

Second, we identified multiple dynamic TFs that are regulators of embryonic neurogenesis but have not yet been studied in adult neurogenesis. DOWN TFs include homeobox protein Dbx2, nuclear glucocorticoid receptor Nr3c1 and Id4. UP TFs included chromatin protein Hmgb1 and proto-oncogene N-myc. During embryonic neurogenesis, these DOWN TFs inhibit cell cycle (Nr3c1) or prevent premature differentiation (Id4), whereas UP TFs regulate progenitor proliferation (Hmgb1, N-myc), suggesting conserved functions during embryonic and adult neurogenesis (Table 2. 2).

Third, more than half of UP TFs and DOWN TFs are largely uncharacterized in the context of neurogenesis, but many of them are close paralogs or binding partners to other neurogenesis-related genes, or have been implicated in other somatic stem cell systems. Examples include

Hmgb1 paralogs (Hmgb2, Hmgb3, Hmga1-rs1), SWI/SNF-related Brg1/Smarca4-associated factors (Smarcc1/Baf155, Smarce1/Baf57), and Nfib paralogs (Nfia, Nfix). In addition, Mxd3, Zeb2 and ZT3/Zfp, regulate adipocyte, melanocyte and myogenic differentiation, respectively, whereas Tsc22d3/Gilz inhibits myogenic differentiation (Table 2. 2). Hopx, which we found to mark qNSCs in the adult dentate gyrus (Figure 2. 5E-G), is expressed in intestinal stem cells and a subset of multipotent hair follicle stem cells (Table 2. 2).

Together, our single-cell transcriptome datasets are a rich resource of genome-wide dynamic expression profiles of TFs during adult neurogenesis. Many TFs that were previously uncharacterized in adult neurogenesis are known to be involved in embryonic neurogenesis or regulation of other somatic stem cells, suggesting shared biology among different stem cell systems and the potential utility of our resource for the general stem cell field.

Molecular Cascades underlying Adult qNSC Activation and Neurogenesis Initiation

The vast majority of UP¹⁰⁰⁰ and DOWN¹⁰⁰⁰ genes in our dataset were not TFs (Table 3. 1). We investigated their characteristics from three perspectives: transition patterns along the developmental trajectory, cellular location of gene products, and biological function.

For transition patterns, plots of the top 150 genes each from UP¹⁰⁰⁰ and DOWN¹⁰⁰⁰ lists showed a wave of molecular activation or inactivation events over time, highlighting the sequential transition of gene expression during qNSC activation and neurogenesis (Figure 3. 1B). To obtain biological insight into these transition patterns, we performed multiple gene ontology analyses. Strikingly, the predicted cellular localizations of protein products of UP¹⁰⁰⁰ genes and DOWN¹⁰⁰⁰ genes were drastically different. 51% of DOWN¹⁰⁰⁰ genes, as opposed to 20% of UP¹⁰⁰⁰ genes, encode proteins associated with the membrane (Figure 3. 3A, p value = 2.6×10^{-29}). On the other hand, 58% of UP¹⁰⁰⁰ genes, as opposed to 20% of DOWN¹⁰⁰⁰ genes, encode proteins associated with the nucleus (Figure 3. 3A, p value = 1.2×10^{-36}). Similar results were obtained using different thresholds for generating lists of UP and DOWN genes (Figure 3. 4A).

Molecular Signatures of Adult qNSCs Revealed by DOWN¹⁰⁰⁰ Genes

Functional annotation of DOWN¹⁰⁰⁰ membrane genes revealed enrichment for ion or protein transport, cell communication and cell adhesion (Figure 3. 3B). Further classification identified proteins specific to the plasma membrane, endoplasmic reticulum, Golgi apparatus, and cytoplasmic vesicles (Figure 3. 4B). KEGG pathway analysis revealed diverse functional entities involved in intra- and inter-cellular communication (Figure 3. 3B). Specifically, Notch signaling, GABAergic synapses, glutamatergic synapses, BMP pathways, MAPK pathway, calcium and cell adhesion-related genes were down-regulated upon qNSC exit from quiescence (Figure 3. 3B, Figure 3. 5A-B and Figure 3. 6A). Electrophysiological recordings of *Nestin-GFP^{cyto+}* NSCs in acute slices from adult animals showed responses to both AMPA and NMDA, suggesting expression of functional receptors (Figure 3. 5B). Each functional signaling pathway entity contained key genes that encode receptors, subunits or downstream mediators (Figure 3. 5A-B). Notably, many ligands for these receptors, including glutamate, GABA, Wnts, BDNF/neurotrophin, Jagged1, BMPs, FGFs and Insulin/IGF2, are known to be present in the adult SGZ niche, suggesting active signaling in qNSCs (Table 2. 2). While previous studies have examined each of these ligands and receptors in regulation of adult neurogenesis in isolation, our systematic genome-wide analyses unified disparate information and suggested a novel model that quiescent adult NSCs are not passive or dormant, but instead are actively integrating various niche signals. More surprisingly, qNSC activation was associated with decreased expression of genes involved in transducing local environment cues and pervasive down-regulation of various signaling pathway-related genes (Figure 3. 5A). These results suggest that, once activated, adult NSCs shunt their capacity to respond to external regulation.

KEGG analysis of DOWN¹⁰⁰⁰ genes also revealed a shift in energy source and metabolism. First, multiple lipid metabolism-related functional entities, including fatty acid degradation and sphingolipid metabolism, were enriched in qNSCs, but down-regulated upon activation (Figure 3. 3B). As previously reported (Knobloch, Braun et al. 2013), qNSCs exhibited the highest level of Spot14 (Figure 3. 6B), which regulates lipid metabolism. qNSCs also maintained an active fatty acid degradation pathway (Acsl3, Acsl6, and Acsbg1; Figure 3. 6B; Table 2. 2). Second, pathway analysis consistently indicated glutathione metabolism and glycolysis as an adult qNSC characteristic, which was lost upon activation. Among glycolysis genes, aldolase A, aldolase C, and

Ldhd decreased significantly, whereas most other glycolysis genes including Gapdh did not change during initiation of neurogenesis (Figure 3. 6C; Table 2. 2).

To validate results from analyses of the top 1000 significantly down-regulated genes, which contained a limited number of genes in each particular pathway, we performed analysis using all expressed genes and an independent functional annotation database wikipathway (Pico, Kelder et al. 2008). Virtually identical results were obtained (Figure 3. 8). Together, analyses of down-regulated genes provided novel insight into molecular signatures of adult qNSCs, including both intrinsic properties and regulation of intra- or inter-cellular signaling pathways.

Sequential Molecular Dynamics during Adult Neurogenesis Revealed by UP¹⁰⁰⁰ Genes

We next analyzed UP¹⁰⁰⁰ genes, which were nucleus-associated and/or related to cell cycle, DNA/RNA metabolism and chromosome organization (Figure 3. 3A). Detailed analysis revealed pervasive activation of cell cycle-related genes, ranging from cell cycle supporting genes, such as nucleotide synthesis, protein/RNA synthesis, and DNA fidelity controls (DNA repair and p53 signaling pathways), to genes directly involved in cell cycle, such as DNA replication, kinetochore complex, cyclin/cyclin-dependent kinases, or cytosolic mitotic spindle genes (Figure 3. 3B).

As opposed to down-regulation of glycolysis-related genes, oxidative phosphorylation-related genes were up-regulated (Figure 3. 6C). Specifically, in contrast to stable expression of earlier mitochondrial respiratory chain complexes (complex I, II, III and IV), expression of subsequent complexes (complex V) increased over pseudotime, implying a gradual completion of the full electron transport chain during neurogenesis (Figure 3. 6C).

The high resolution of Waterfall analyses revealed temporal relationships among genes in different functional groups. Cell cycle checkpoint genes were sequentially activated following the known biological sequence of cell cycles: G1 to S transition, followed by G2 to M transition and then chromosomal segregation, indicating that our pseudotime accurately reconstructs sequential biological events (Figure 3. 6D). Initiation of cell cycle preceded the major transcriptomic shift (Figure 3. 5C). Notably, up-regulation of genes encoding ribosomal subunits preceded the appearance of any cell cycle checkpoint genes (Figure 3. 5C and Figure 3. 6D), suggesting that

priming of protein synthesis machinery may mark the G0 to G1 transition ahead of cell cycle entry during adult qNSC activation.

Together, analyses of UP¹⁰⁰⁰ genes suggested that molecular dynamics of qNSC activation and initiation of neurogenesis are largely defined by priming of protein synthesis machinery, cell cycle entry, activation of RNA and protein biogenesis, and a shift in energy metabolism from glycolysis to oxidative phosphorylation. An independent approach using a different functional annotation database showed similar results (Figure 3. 8).

Holistic Picture of Molecular Cascades underlying Adult Neurogenesis

Initiation

Based on the molecular dynamics from qNSCs to aNSCs and then eIPCs, we have reconstructed sequential waves of biological events from single-cell RNA-seq data and Waterfall (Figure 3. 7). The process begins with adult qNSCs down-regulating transcription factors defining quiescence and decreasing competence for cell signaling (RTKs, GPCRs, neurotransmitter receptors, cytokines, calcium). Concurrently, glycolysis, glutathione and fatty acid metabolism begins to wane, while up-regulation of protein translation capacity is the first marker of a pre-activation stage. As NSCs enter cell cycle, oxidative phosphorylation becomes the primary energy source. Progression through cell cycle accompanies a major decline in NSC metabolism (glutathione, fatty acid, drug metabolism) and an increase in eIPC transcription factors. Finally, kinetochore and chromosomal segregation occurs in the first neurogenic progeny. Overall, the developmental trajectory is defined by a coordinated switch from a membrane-targeted to a nuclear-targeted transcriptome, suggesting a transition from qNSCs dominated by extrinsic signaling to eIPCs dominated by a pre-programmed intrinsic molecular cascade.

Downregulated genes

CST2	0.3797731	Prox2	0.361586269	Gli3	0.296752247	Atf9pD6	0.257399688	Gwcm2	0.22442398	Mmd	0.20099996	Tarfb	0.184695989	Uppr2	0.168712039	Enof	0.154746029	Tmem184	0.14283671
Tspan7	0.78203821	Atf9pA5	0.37955954	Ntn	0.29595996	Gq3	0.2564988	Mu2	0.22449829	Acan	0.20022935	Anu1	0.18467026	Xnn	0.18387563	Bcl2L2	0.14267266	Tmem3	0.14269382
Sc1a1	0.7480762	Sc1a5	0.37977915	Birc2	0.29465366	Au17	0.23860968	Utrn	0.22384646	Utrn	0.20040761	Prox1b	0.18469283	Utrn	0.18387563	Sc1a1b	0.14267266	Tmem3	0.14269382
Sc1a1	0.71892922	A63008078a	0.37789399	Cxcr4	0.29389876	Cxcr4	0.25722362	Nr3a4	0.22385408	Badrh	0.19973064	Tmem15	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.7140141	Dna3	0.37750099	Gpr2	0.29343245	Prox2	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Quant	0.70558566	Tmem2	0.37525996	Muc	0.29323289	Mu2	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Gd3a3	0.70709699	Kat5	0.37400884	Au17	0.29323289	Tspan7	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Atf9pA5	0.69173774	Au17	0.373541788	D83003920a	0.29323289	D83003920a	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Thy1	0.6844683	Nr3a4	0.37242099	M1	0.2895036	A63007789a	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Atf9pA5	0.68722721	Leu1a	0.3699746	Cxcr4	0.29389876	Prox2	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Apoe	0.67959598	G2D2	0.36858683	Birc2	0.29389876	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Cu1	0.66036808	Atf9pA5	0.37242099	M1	0.2895036	A63007789a	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
In2b	0.65579772	Lamp	0.367655364	Kdm3a	0.28776893	43934024078a	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Meot1	0.64272427	Sc1a	0.36778791	Sc1a1	0.28776893	E33007048a	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Nr3a4	0.63945884	Dna1	0.36646665	Hprt2c	0.28776893	A39180	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Pp1a1	0.63549348	Cxcr4	0.36520231	Gpr2	0.29343245	Muc	0.29323289	Mu2	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266
Nr3a4	0.63549348	Dna1	0.36520231	Gpr2	0.29343245	Muc	0.29323289	Mu2	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266
Gpr1	0.58905954	Lgr4	0.36466248	Enr1	0.28864431	Enr1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Gpr2	0.58942444	Nr3a4	0.36466248	Enr1	0.28864431	Enr1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1	0.28776893	Sc1a1	0.25722362	Nr3a4	0.22385408	Nr3a4	0.19973064	Prox2	0.1843949	Rarb	0.18352424	G2D4	0.14267266	Aut	0.14165435
Prox2b	0.57863669	Cxcr4	0.36520231	Sc1a1															

Table 3. 2. List of UP TFs and DOWN TFs.

UP TFs	DOWN TFs
Hmgb2	Id4
Hmgb1	Hopx
Sox11	Sox9
Ybx1	Npas3
Nfib	Bhlhe41
Nfix	Zfp71rs1
Insm1	Id3
Hmgb3	Hes1
Ssrp1	Irf3
Smarcc1	Tsc22d3
Zfp367	Rfx4
Smarce1	Nr2e1
Mxd3	Tcf7l2
Zfp62	Pou3f3
Zeb2	Zscan26
Nfia	Dbx2
Sox4	Sox2
Terf1	Nr3c1
Tcf4	Zbtb4
Tulp4	Tfe3
Emx1	Dmtf1
Hmga1-rs1	Bcl6
Zfp637	Fos
Mycn	Ikzf2
Foxg1	Nfx1
Tcf12	Zfp740
Cebpg	Grhl1
Eomes	Klf15
Zmiz1	Zbtb26
Nfyc	4932411N23Rik
Zfp386	Arx
Zkscan1	Etv5
Bcl11a	Mier1
Zfp763	Fezf2
Zfp944	Zfpm2
Zbtb41	Zfp869
Trp53	Lef1
Csde1	Thap7
Zmiz2	Mxi1
Zfp808	Camta2
Zbtb38	Zfp617
Atf4	

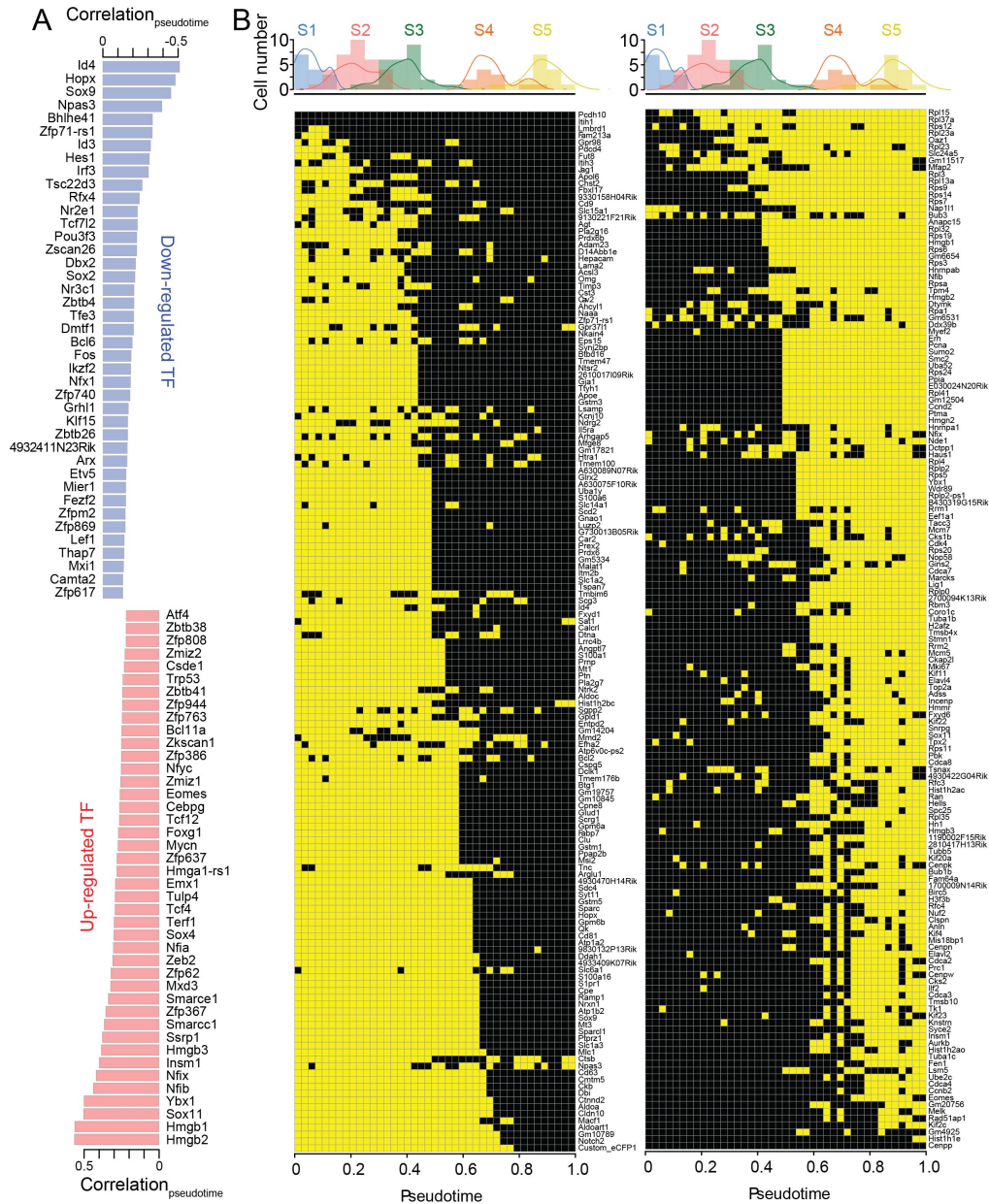


Figure 3. 1. Molecular cascade underlying adult quiescent neural stem cell activation and neurogenesis initiation.

(A) Lists of DOWN and UP TFs and their Spearman correlation coefficient with pseudotime.

(B) ON/HIGH (yellow) or OFF/LOW (black) states of top 150 DOWN (left) and UP (right) genes sorted by the timing of transition points. Shown on the top are histograms of the numbers of individual cells examined along the pseudotime progression. The colors on the histogram follow the color scheme of Figure 2. 3.

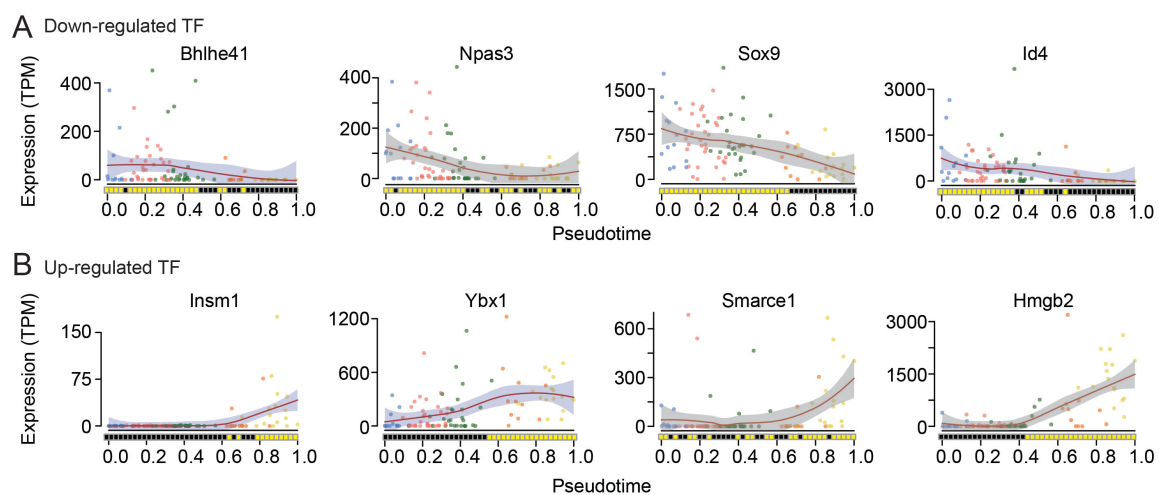


Figure 3. 2. Pseudotime profiles of sample transcription factors.

Shown are pseudotime profiles of representative transcription factors that were down-regulated (A) or up-regulated (B) over developmental pseudotime progression.

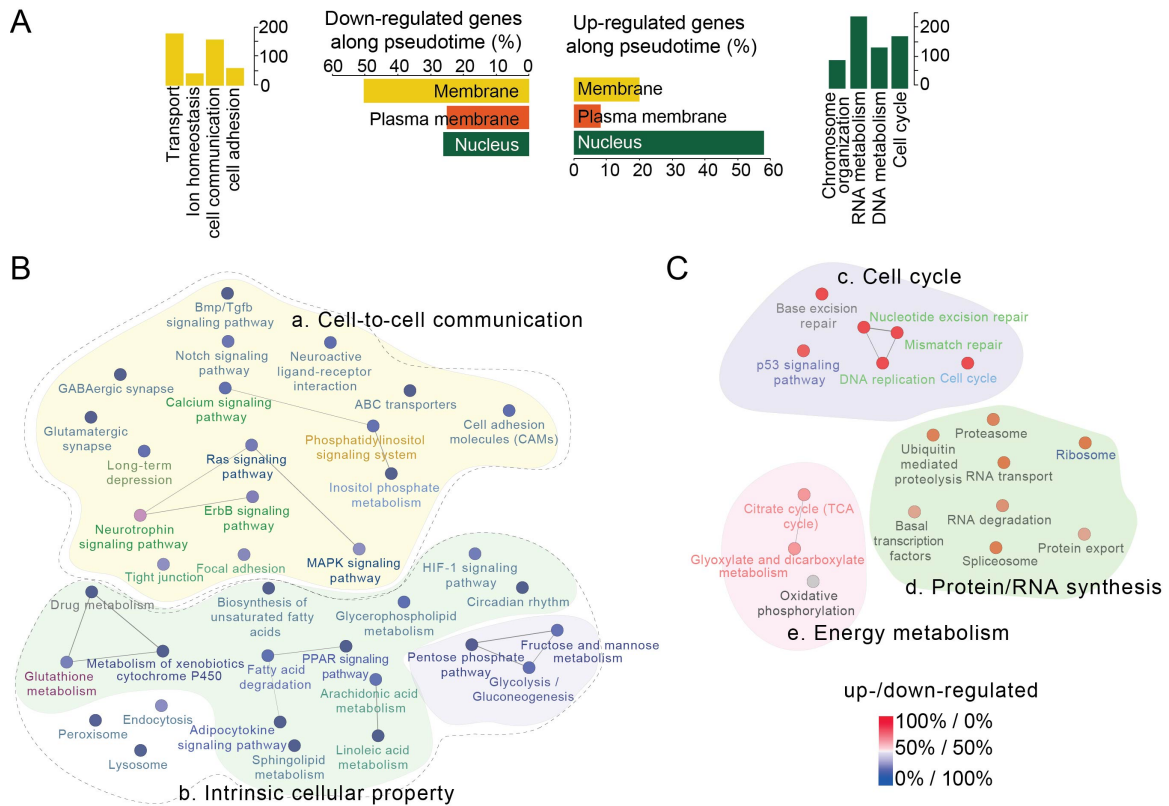


Figure 3. 3. Functional characterization of UP1000 genes and DOWN1000 genes.

(A) Quantification of predicted cellular location of gene products of UP¹⁰⁰⁰ genes and DOWN¹⁰⁰⁰ genes (two middle panels). Also shown are numbers of genes with indicated functions for membrane-associated DOWN genes (left panel) and those for nucleus-associated UP genes (right panel).

(B) Functional GO analysis for DOWN¹⁰⁰⁰ and UP¹⁰⁰⁰ genes. Color of each functional entity represents proportion of UP genes (blue) and DOWN genes (red). Connections between each pair of data points represent sharing more than 5 genes between the pair. Functionally similar entities are grouped with same background colors. Broken lines represent two categories of DOWN genes: genes encoding proteins involved in the intra- or extra-cellular communication and genes encoding proteins defining intrinsic stem cell properties. The P values are from hypergeometric tests, and corrected by Holm–Bonferroni method.

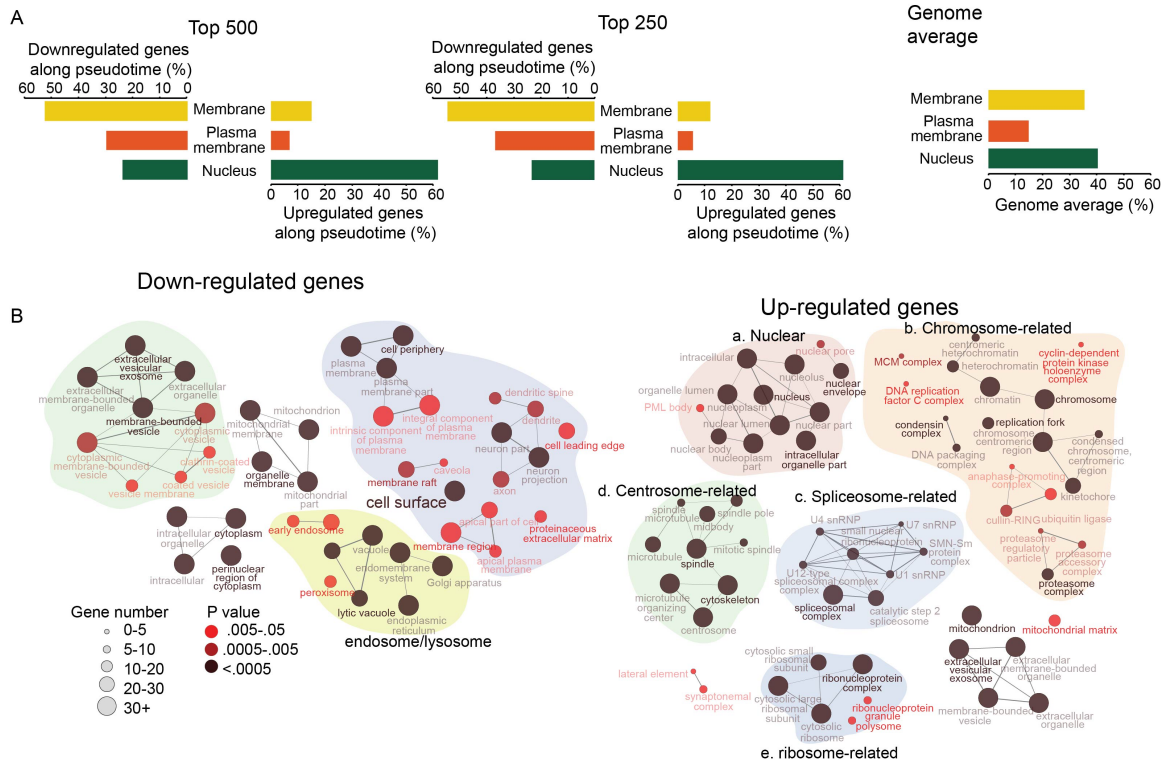


Figure 3. 4. Cellular compartment analysis for UP and DOWN genes. (A) Cellular compartment analysis for UP and DOWN genes with different thresholds. Shown are summaries of enrichment patterns of UP and DOWN genes with cutoff for top 500 genes (left panel) or top 250 genes (middle panel), and average for all genes (right panel).

(B-C) Detailed predicted cellular locations for gene products of DOWN¹⁰⁰⁰ genes (B) and of UP¹⁰⁰⁰ genes (C). The size of each data point represents the number of genes within each predicted location, and color of each data point represents P value of enrichment to each predicted location. P values are from hypergeometric test and corrected by Holm–Bonferroni method (Bonferroni step-down correction). Connections between pairs of data points represent sharing more than 50% of genes between the pair. Similar entities are grouped with the same background colors.

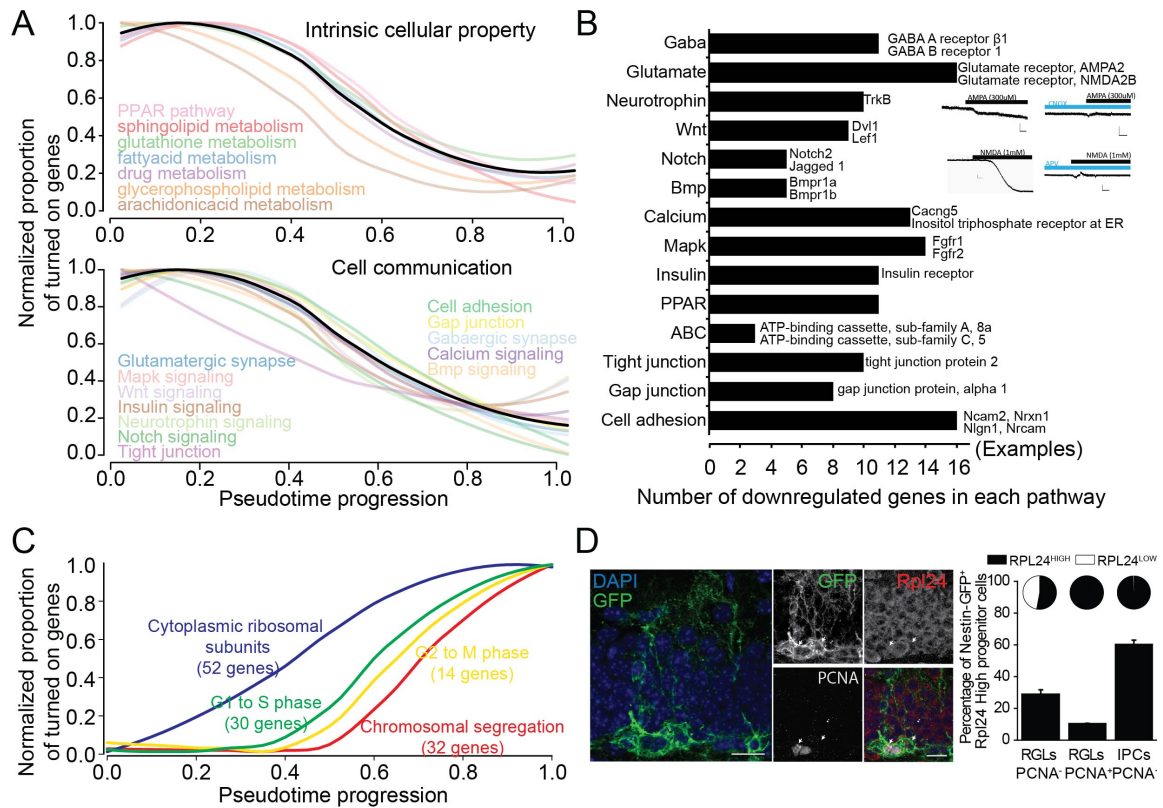


Figure 3.5. Sequential molecular dynamics during adult neural stem cell activation and neurogenesis.

(A) Gradual down-regulation of averaged binary states of each functional entity defining intrinsic stem cell properties (top) and those defining intra- or extracellular communication (bottom). On/off binary states of expressed genes within each functional entity were averaged and normalized to show the timing of the transition.

(B) Number of genes in different DOWN gene ontology groups and representative example genes. Expression patterns of representative genes over developmental pseudotime are shown in Figure 3.6A. Electrophysiological recording of Nestin-GFP^{cyt} NSCs in acute hippocampal slices showed responses to both AMPA and NMDA (bottom).

(C) Gradual up-regulation of averaged binary states of each functional entity defining cell cycle checkpoints and cytoplasmic ribosomal subunits. On/off binary states of up-regulated genes within each functional entity are averaged and normalized to exemplify the timing of the transition. Cell cycle checkpoint genes were up-regulated in the sequence of the cell cycle progression. The up-regulation of cytoplasmic ribosomal subunits preceded up-regulation of the earliest cell cycle checkpoint gene.

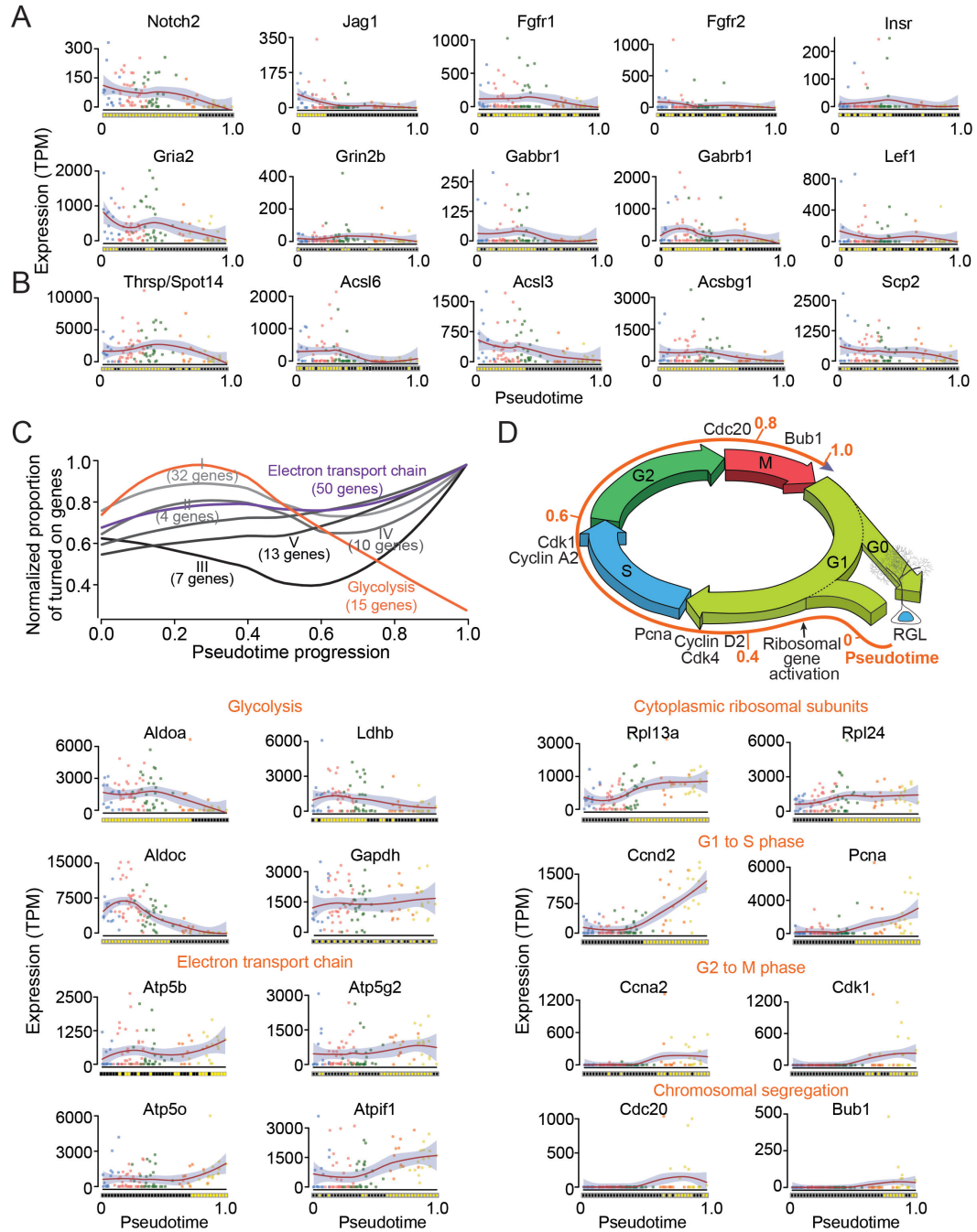


Figure 3.6. Representative pseudotime profiles of key genes related to niche signaling, cell cycle progression, and energy metabolism.

(A-B) Pseudotime profiles of representative genes related to niche signaling (A) and genes related to fatty acid metabolism (B), both of which exhibited down-regulation upon qNSC activation. Similarly plotted as in Figure 2. 3D.

(C) Reciprocal changes of expression of mitochondrial electron transport genes and glycolysis genes. Roman numerals represent electron transport complex. Complex V exhibited a clear increase throughout early neurogenic pseudotime. Shown at the bottom are pseudotime profiles of some genes related to glycolysis (gradual down-regulation) and complex V genes in mitochondrial electron transport chain (gradual up-regulation).

(D) Developmental pseudotime recapitulates cell cycle progression during qNSC activation and initiation of neurogenesis. Shown on the top is a schematic illustration of cell cycle phases correlated with pseudotime. Cell cycle checkpoint genes were up-regulated following the sequence of biological cell cycle progression. Shown at the bottom are pseudotime profiles of genes related to cell cycle checkpoints and cytoplasmic ribosomal subunits. Note that ribosomal genes were up-regulated earlier than any of the cell cycle related genes.

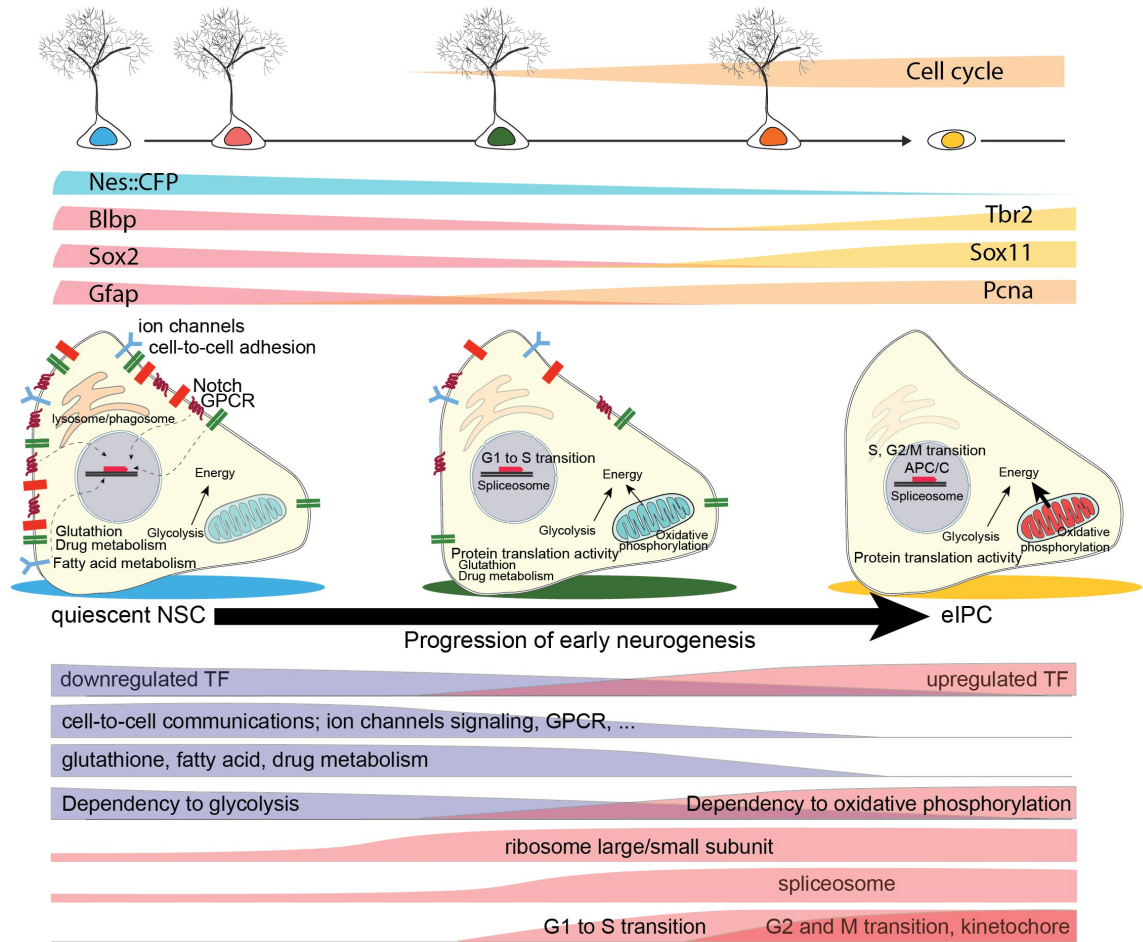


Figure 3. 7. Schematic summary of molecular signatures of quiescent adult neural stem cells and molecular cascades underlying their activation and neurogenesis.

Shown on top is an illustration of known marker expression and cell cycle activation during adult hippocampal neurogenesis. Shown in the middle is an illustration of molecular signatures of adult qNSCs and their immediate progeny. Shown at the bottom are functional categories of genes that show a clear shift during adult qNSC activation and generation of eIPCs. qNSCs exhibit intra- and inter-cellular signaling to actively sense the local niche, rely mostly on glycolysis for energy, and have highly active fatty acid, glutathione, and drug metabolism. Upon activation, NSCs increase translational capacity, followed by cell cycle entry with G1 to S transition. Oxidative phosphorylation starts to be active and stem cell specific properties are down-regulated. eIPCs maintain active cell cycle genes, ribosomal activity and fully active oxidative phosphorylation for energy generation. The color scheme on the top and the middle illustration follows the colors from Figure 2. 3B.

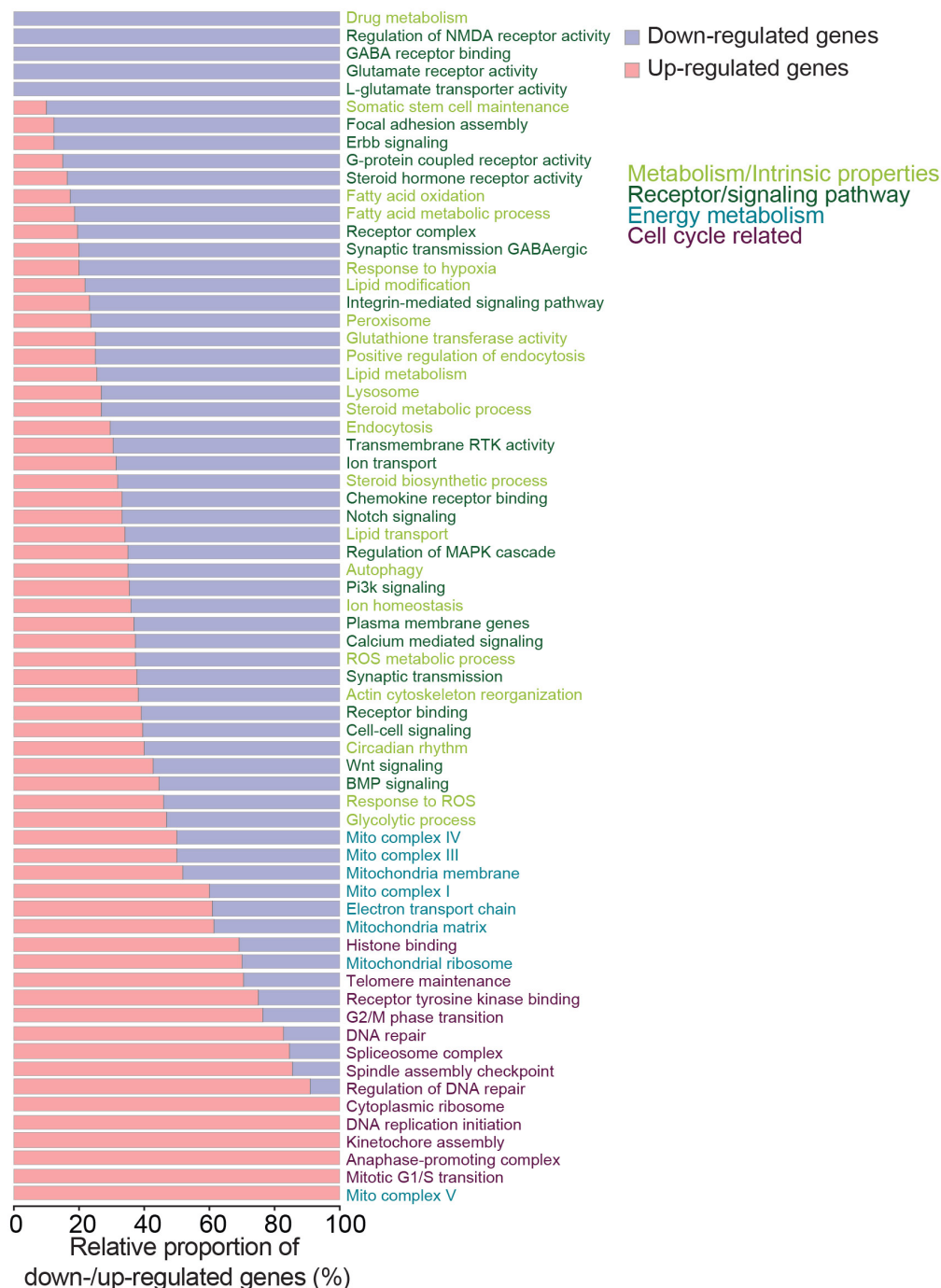


Figure 3. 8. Independent validation for GO entity enrichment test for up-regulated and down-regulated genes during quiescent stem cell activation and neurogenesis.

Shown is a summary of the proportion of up-regulated and down-regulated genes within each key functional GO entity. GO entities with a disproportionately higher proportion of up-regulated genes represent functional pathways that are activated during exit of quiescence and early stages of neurogenesis, whereas GO entities with a disproportionately higher proportion of down-regulated genes represent functional pathways that are qNSC-specific pathways.

Chapter 4. Discussion⁴

⁴ This chapter is based on Shin, J., et al. (2015). "Single-cell RNA-seq with waterfall reveals molecular cascades underlying adult neurogenesis." Cell Stem Cell **17**(3): 360-372.

Understanding adult NSC behavior and neurogenesis requires quantification of molecular states along a continuous developmental process. In the current study we generated three major resources. First, we provide a comprehensive dataset of single-cell transcriptomes of qNSCs and their immediate progeny in the adult mouse hippocampus *in vivo*. Second, we provide Waterfall, an unsupervised bioinformatic suite for *in silico* reconstruction of molecular trajectories based on snapshots of single-cell transcriptomes and statistical gene expression analysis over continuous developmental processes. Third, we provide a holistic picture of adult qNSC molecular signatures and dynamic molecular cascades underlying initial phases of adult neurogenesis at unprecedented temporal resolution (Figure 3. 7). Our study provides an example of how to resolve cellular heterogeneity and reveal developmental dynamics for systematic molecular characterization of stem cells and their differentiation *in vivo*. Our approach can be adapted for various single-cell omics analyses (transcriptomics, proteomics, epigenomics, lipidomics, and metabolomics) of many continuous biological processes, such as development, physiological and pharmacological stimulation, and disease progression (See examples in Chapter 5. Single cell analysis with Waterfall).

A Resource of *in vivo* Single-cell Transcriptomes of qNSCs and their Immediate Progeny

Our study provides a single-cell RNA-seq dataset and the first comprehensive view of transcriptome dynamics underlying adult qNSC behavior *in vivo*. Currently, there is no published dataset for transcriptome dynamics during stem cell development in any somatic system *in vivo*. We performed multiple levels of validation of our dataset and approach, including comparison with an *in situ* database, confirmation with known and unknown marker expression during adult neurogenesis *in vivo*, and functional validation via clonal lineage tracing and electrophysiology (Figures 3 and 6B; Table 2. 2).

Our resource of holistic molecular profiles during early neurogenesis has three unique features that increase its versatility. First, our whole-transcriptome information includes unannotated transcripts, isoforms, and retrotransposon-derived transcripts, as opposed to multiplexed qPCR or microarray-based studies which can only provide limited annotated transcripts (Hoppe, Coutu et al. 2014). Second, we animated static single-cell transcriptomes over

the in vivo neurogenesis trajectory, allowing queries of molecular dynamics for each gene over the continuous biological process and generation of novel functional hypotheses during specific phases of adult qNSC maintenance, activation and neurogenesis initiation. For example, expression of nuclear glucocorticoid receptor *Nr3c1* in adult qNSCs but not in eIPCs (Figure 3. 1A) suggests a cellular target of glucocorticoids during adult hippocampal neurogenesis and a means to manipulate qNSCs in vivo. Our resource also reveals potential prospective markers of adult neurogenesis, such as using *Hopx-CreER^{T2}* to target adult qNSCs in vivo (Figures 3E-G). Third, each transcriptome in our dataset represents a biological state within a single cell. This modular construction allows for flexible reorganization of the dataset to probe different questions as our understanding of the biology evolves. For example, instead of generating reporter lines for individual genes or using different surface markers for physical sorting of specific cell populations, investigators can perform unlimited in silico cell sorting with any individual gene, or multiple genes in combination, to obtain a selected cell population to probe their gene expression characteristics at the genome-wide level using our single-cell datasets.

Developmental Dynamics of Adult Neurogenesis at the System Level

Recent genome-wide studies have begun to provide a system-level understanding of in vivo adult NSC biology using marker-defined NPC populations (Kriegstein and Alvarez-Buylla 2009, Bracko, Singer et al. 2012, Codega, Silva-Vargas et al. 2014). Yet previous studies have only acquired snapshots of transcriptomes, which limits investigation of developmental dynamics among different cellular states. We co-opted the imperfection of the *Nestin-CFP^{huc}* genetic labeling system to collect individual qNSCs and their immediate neuronal progeny concurrently. Aligning cells along the developmental trajectory yielded, for the first time, a molecular continuum with sequential progression of the individual transcriptome from qNSC to aNSC and then eIPC. Importantly, this novel approach does not divide developmental processes into discrete stages that are defined a priori by capturing populations sharing specific markers.

Our resources provide unparalleled temporal resolution to identify new mechanisms underlying adult NSC biology. For example, we showed that Acyl-CoA synthetases (*Acs13*, *Acs16* and *Acsbg1*), the enzymes for the first step of fatty acid β -oxidation, were highly expressed only in qNSCs (Figure 3. 6B), suggesting a novel role for active fatty acid β -oxidation in qNSCs and

thereby extending previous findings on the role of fatty acid metabolism in adult neurogenesis (Knobloch, Braun et al. 2013)(Table 2. 2). We also found that ribosomal subunits were the first genes up-regulated upon adult qNSC activation and early differentiation (Figure 3. 5C and Figure 3. 6D), suggesting a possible demarcation of G0 to G1 transition and providing the timing for switches in protein synthesis, the regulation of which is important for somatic stem cell function (Signer, Magee et al. 2014). The holistic picture we obtained unifies disparate information and illuminates novel biological themes in stem cell biology. For example, our analysis suggests that qNSCs actively respond to local environmental cues through various signaling pathways, but gradually and globally shut off signaling capacity upon activation. These observations support the concept of a niche wherein mammalian somatic stem cells are tightly controlled by a regulatory microenvironment (Schofield 1978), and predict that ePCs are less responsive to environmental input (Berg, Yoon et al. 2015). This novel biological insight may be applicable to many somatic stem systems defined by stochastic behavior (Simons and Clevers 2011).

Waterfall Analysis of Single-cell Transcriptomes within a Continuum

Waterfall has three key differences from previously methodologies. First, it requires very little prior information to generate a highly accurate temporal trajectory at single-cell resolution. Previous methods have been able to reconstruct accurate trajectories by relying on a robust set of known markers to establish cell order and validate cell alignment at numerous points along the timeline. For many biological systems and processes, we have much less information. Second, in contrast to Monocle (Trapnell, Cacchiarelli et al. 2014) or Wanderlust (Bendall, Davis et al. 2014), Waterfall uses k-means clustering to build a trajectory and assign an individual cell a pseudotime based on each cell's proximity to the cluster-derived trajectory, rather than constructing a trajectory by directly connecting each cell to the next. Third, in order to analyze stochastic gene expressions, we adopted HMM to predict consecutive binary states in gene expression activity over pseudotime. HMM permits the interpretation of highly variable data without logarithmic transformation, normalization, or the input of any arbitrary parameters, such as threshold for gene expression noise or Markovian parameters (transition probability, initial probability and emission probability). Our validation for known and unknown genes in adult neurogenesis indicated that HMM correctly predicted temporal dynamics of in vivo biology.

There is no conceptual restriction of our approach to transcriptome studies of adult neurogenesis. Indeed, in Chapter 5. Single cell analysis with Waterfall, we provide examples of how Waterfall could be broadly applicable for single-cell RNA seq datasets such as in vitro myogenesis, in vivo embryonic lung development, single-cell mass-cytometry dataset from in vivo B cell development, and synthetic datasets. We expect that Waterfall algorithms can be adopted for diverse single-cell multi-dimensional datasets, including single-cell transcriptomes, epigenomes, proteomes, and metabolomes, of various continuous biological processes.

Chapter 5. Single cell analysis with Waterfall⁵

⁵ This chapter is based on Shin, J., et al. (2015). "Single-cell RNA-seq with waterfall reveals molecular cascades underlying adult neurogenesis." Cell Stem Cell **17**(3): 360-372.

I. Development of Waterfall

Investigating the molecular signatures of adult hippocampal neurogenesis at the single cell level presents considerable challenges – challenges both endemic to the field of single cell analysis and unique to the adult neurogenesis system. In this chapter, we briefly highlight the most critical of these challenges and our strategy to resolve them. We then provide a detailed summary of our development and application of Waterfall to our adult neurogenesis dataset, followed by a comparison of the results with those obtained by Monocle for the same dataset. We then evaluate how well our bioinformatic strategy handles data from different biological systems, each with its own characteristics and constraints. The datasets we analyzed are shown in Table 5. 1. List of datasets to which Waterfall was applied. .

What is Waterfall?

Waterfall is a bioinformatics pipeline to process multi-dimensional single cell dataset derived from continuous biological processes. Waterfall comprises three steps (Figure 5. 1):

Pre-processing

- a. Removing outliers
- b. Determining and orienting a single route of interest and constitutive single cell clusters

Reconstruction

- a. Reconstruct a trajectory using minimum spanning tree (MST)
- b. Determining the relative chronological location (pseudotime) of each single cell

Gene expression analyses

- a. List of positively or negatively correlated genes with pseudotime
- b. Determining underlying binary gene expression states using hidden Markov model

Challenges of analyzing single cells in continuous biological processes

Using single cell molecular profiles, we can deconstruct heterogeneity at the molecular level with greater precision than is possible from analyzing bulk populations. The goal and methodology of single cell analysis is tightly bound to the type of heterogeneity that to be resolved.

Resolving static or cell-type heterogeneity is analogous to directly comparing group data at the population level, and can be performed with conventional analytical tools. The goal is to identify distinct categories of cells based on the molecular profiles and classify cells accordingly. Recently developed cell-type prediction software has great potential to facilitate these studies (Cahan, Li et al. 2014). In contrast, there are no tools for solving dynamic or temporal heterogeneity that can be adopted from population level studies, because it was not until the advent of single cell technology that we had access to molecular profiles at high enough resolution to track gradual transitions along continuous biological processes. Single cell analyses of dynamic processes face two critical challenges: (1) lack of information and (2) stochasticity of the single cell transcriptome (Figure 5. 2).

Lack of information

In order to deconstruct heterogeneity within a continuous process like development, we need to find genes that change during the process to provide points of reference and alignment (Buettner, Natarajan et al. 2015).

Bendall et al. used known set of markers involved in the biological process, thereby successfully de-convoluting the in vivo B cell developmental process from single cell mass cytometry data (Bendall, Davis et al. 2014) (Figure 5. 2A). All the single cell mass cytometry and single cell qPCR studies use this approach (Bendall, Simonds et al. 2011, Qiu, Simonds et al. 2011, Bodenmiller, Zunder et al. 2012, Guo, Luc et al. 2013, Bendall, Davis et al. 2014, Saadatpour, Guo et al. 2014). Applying this methodology, however, requires a substantial understanding of the biological system, where we know most, if not all, of the important markers to define cell types or developmental stages. For most biological system other than the hematopoietic system, we do not yet have enough information to identify key genes or stratifiers.

Trapnell et al. nicely solved this challenge by discovering key genes through population comparison (Figure 5. 2B) (Trapnell, Cacchiarelli et al. 2014). Trapnell et al. astutely and successfully reconstructed the trajectory of in vitro skeletal muscle differentiation at single cell resolution using differentially expressed genes from population comparison between before and after switching media to initiate differentiation (Figure 5. 2B). To perform this kind of population comparison, however, we need to know the temporal boundaries that define subgroups before analyzing at the single cell level. Such information is a clear advantage of in vitro systems and rarely available for in vivo systems, considering the heterogeneity of cell types and asynchronous cellular states (Figure 5. 2C).

Stochasticity of single cell transcriptome

Single cell gene expression is highly stochastic, whereas the population level transcriptome is more stable due to averaging thousands to millions of single cell transcriptomes (Ramskold, Luo et al. 2012, Shalek, Satija et al. 2013). One of the most widely used models to explain single cell level gene expression is a bistable transcription model, wherein underlying gene expression activity oscillates between an active state and a basal state (Novick and Weiner 1957, Levsky, Shenoy et al. 2002, Ozbudak, Thattai et al. 2004, Lim and van Oudenaarden 2007, Arnoldini, Vizcarra et al. 2014). Assuming the mRNA transcription is a major determinant for mRNA abundance (Schwanhauser, Busse et al. 2011), we could approximate single cell transcriptome stochasticity as the random variability from two disparate underlying gene expression states. When analyzing single cell datasets from continuous biological processes, we need to furthermore consider the transition events between the underlying states during the process. Thus, we need a novel approach that can better reflect the single-cell biology. An ideal approach should explain:

1. Transition events between the two underlying states throughout the continuous process
2. Different points of the transition events of individual genes over the continuous process
3. Different gene expression distributions resulting from the two underlying states
4. Noisy or missing data points

We proposed to use a hidden Markov model (HMM) when interpreting single cell data from continuous processes (Figure 5. 3). Major assumptions of using HMM for single cell data from continuous trajectory are:

1. There are two underlying states emitting different single cell gene expression distributions.
2. The bi-stable state model is the major sources of single cell gene expression stochasticity.
3. The pseudotime bin (Figure 2. 6) is small enough to contain cells with identical “hidden” states for most genes.
4. Gene expression from each state follows a normal distribution.

We believe the first two assumptions are reasonable, as shown by multiple single cell studies (Novick and Weiner 1957, Levsky, Shenoy et al. 2002, Ozbudak, Thattai et al. 2004, Lim and van Oudenaarden 2007, Arnoldini, Vizcarra et al. 2014). The third assumption depends on the size of bins, thus can be iteratively modified through comparison between the prediction and the known biology. The fourth assumption is also reasonable, since we are not forcing the gene expression distribution to a Poisson distribution. In addition, HMM is relatively immune to noisy or missing data points.

Admittedly, HMM might be one of many potential models to explain stochastic single cell RNA-seq data over continuous trajectories. However, there is no attempt yet to try to systematically interpret stochasticity in a single cell RNA-seq dataset. We introduced this concept for the first time and it fits well with the known biology as could be validated for specific genes.

II. Analyzing single cell data using Waterfall

We originally developed the Waterfall pipeline to analyze the single cell RNA-sequencing dataset from in vivo adult hippocampal neurogenesis. In this chapter, we present detailed analytical procedures of how we applied Waterfall for the adult neurogenesis system. Then, we tested how Waterfall performed in other systems (e.g. in vitro myogenesis, embryonic lung development, in vivo B cell development) and platforms (e.g. single cell mass cytometry) (Bendall, Davis et al. 2014, Trapnell, Cacchiarelli et al. 2014, Treutlein, Brownfield et al. 2014). In addition, we compared the performance of Waterfall and Monocle, which is a pioneering and currently the only available software for single cell RNA-seq from continuous processes (Trapnell, Cacchiarelli et al. 2014). Finally, we also applied Waterfall on two synthetic datasets to show its optimal behavior under ideal conditions.

1. Adult neurogenesis dataset (In vivo single cell RNA-seq)

Nestin-CFP^{nuc} is highly expressed in the neural stem cell (NSC) population and is carried over to their immediate progeny (early immediate progenitor cells, eIPCs). Using this imperfect labeling system to our advantage, we were able to capture the continuous cellular states between quiescent NSCs (qNSCs) and eIPCs.

Mapping and calculating normalized expression levels

The CFP^{nuc+} cells expressed higher levels of CFP transcript compared to CFP^{nuc-} cells (Wilcoxon signed-rank test p value = 8.018×10^{-16}) (Figure 5. 4A). The 3' bias of the amplification was comparable to other single cell RNA-seq studies (Figure 5. 4B). Of note, DNase treatment significantly improved the data quality (Figure 5. 4B). We used RSEM to calculate TPM values (Li and Dewey 2011).

Waterfall: Preprocessing

Nestin-CFP^{nuc} mouse system also labels a small percentage of non-NPCs in the adult dentate gyrus (Figure 2. 2A). We excluded cells that either exhibited markedly different transcriptomic profiles from majority of the CFP^{nuc+} population or were identified as non-NPCs, such as oligodendrocyte progenitor cells or pericytes (Figure 5. 5).

We then performed unsupervised clustering and principal component analysis (PCA; Figure 5. 6). S1 to S5 formed a linear trajectory, whereas SA formed a tangential branch (Figure 5. 6B) as also indicated by the unbiased minimum spanning tree (MST; Figure 5. 7A). Gene expression profiles at Figure 5. 6C indicated that the S1-S2 population represents the quiescent neural stem cell (qNSC) population, whereas the S5 branch represents the early intermediate progenitor cell (eIPC) population, which gives rise to newborn neurons in adult dentate gyrus.

Waterfall: Additional trajectories

Although S1 to S5 path represents the most probable path of adult neurogenesis, there are other possible routes to consider (Figure 5. 7). First, we have three hypotheses for the identity of SA (Figure 5. 7C, left and middle): (1) astrocytic differentiation from NSC depletion (Encinas, Michurina et al. 2011)(Figure 5. 7C, left); (2) activated NSCs returning to quiescence after asymmetric or symmetric division (Bonaguidi et al., 2011); (3) a distinct adult NSC population (Figure 5. 7C, middle). SA expressed adult qNSC genes, such as Gfap at similar or even higher level than S1 and S2. They also expressed significantly higher levels of oxidative phosphorylation components (Cyb5, Atp5o; unadjusted p-values 0.0275 and 0.00211, respectively) and mitochondrial ribosomal subunits (Mrps11, Mrps21; unadjusted p-values 0.0156 and 0.0234, respectively), suggesting that SA was at a distinct metabolic state.

Second, while building up a pseudotime sequence, we initially noticed that a few cells located at the beginning of the trajectory expressed lower levels of adult NSC markers (Figure 5. 7C, right). Thus, we excluded the five cells, which were located at the very end of S1 for pseudotime analyses. Our working hypothesis is that they might be exhaustive gliogenic population (Bonaguidi, Wheeler et al. 2011, Encinas, Michurina et al. 2011), since these five cells expressed high level of glial markers. We do not have resolution to draw a line between qNSC populations from potential glial differentiating population. Although these possibilities propose exciting possible directions for more data mining from our single-cell RNA-seq dataset, we decided to focus on the neurogenic trajectory.

Waterfall: Reconstructing trajectory and pseudotime

We used k-means clustering and minimum spanning tree to generate a trajectory (Figure 5. 8A). The number of clusters was determined iteratively using the elbow method (Thorndike 1953). We assumed that the tangential components of single cells to the trajectory were derived from the stochastic gene expression or technical variability. We assigned pseudotime to each single cell based on its relative location over the pseudotime trajectory (Figure 5. 8B).

Waterfall: Gene expression analyses

Pseudotime allows profiling of gene expression dynamics throughout the reconstructed trajectory (Figure 5. 9). As expected, the single cell transcriptome was highly stochastic (Figure 5. 9). The regression plot was not quantitative enough because the shape and deflection point of the regression plot varies greatly based on arbitrary parameters such as degree of smoothing. Moreover, the regression plot does not reflect the stochastic and bimodal nature of single cell expression. HMM is more quantitative and does not require setting any arbitrary parameters (Figure 5. 9, yellow/black block graphs).

Applying Monocle to the adult neurogenesis dataset

We applied Monocle to our adult neurogenesis dataset to compare its performance with Waterfall (Trapnell, Cacchiarelli et al. 2014). In order to perform a fair comparison, we used TPM both for Monocle and Waterfall, because TPM was shown to be more stable in analyzing RNA-seq data (Wagner, Kin et al. 2012). We then removed oligodendrocytes, pericytes or other unidentified outliers based on their marker expression before applying Monocle (Figure 5. 5).

In vivo derived single cell data do not entail spatiotemporal information necessary to select for key developmental genes to stratify the single cells. Since Trapnell et al. recommends selecting highly expressed genes, we chose a trajectory with top 1000 genes (Figure 5. 10).

We compared the gradual transition of single cell transcriptome along the trajectories determined by Waterfall and Monocle, which was evaluated by the sum of the distance between each pair of neighboring cells from the beginning to the end of the pseudotime sequence (Figure 5. 11A). The more gradual the trajectory, the shorter the total distance would be. As expected, Waterfall and Monocle resulted in significantly shorter distances compared to a randomly permuted cell order (Figure 5. 11B). Further, the total path distance of the pseudotime by Waterfall was significantly smaller than the pseudotime by Monocle. The cumulative probability in Figure 5. 11B demonstrates that it requires $\sim 10^{18}$ trials of random permutation of single cell orderings to achieve the precision of Waterfall and $\sim 10^{10}$ trials of random permutation of the single cell orderings to achieve the precision of Monocle.

Monocle showed expected profiles of qNSC-enriched genes (Apoe, Fabp7, and Gfap) and eIPC-enriched genes (Tbr2/Eomes and Sox11) (Figure 5. 12). However, Monocle misplaced a few cells, resulting in bimodal expression pattern of Tbr2/Eomes and Sox11 over time, which is counter to the known biology (Figure 5. 12A). In contrast, Waterfall not only placed these cells correctly but also provided on/high and off/low states for a better appreciation of the gene expression state transition (Figure 5. 12B).

2. In vitro myoblast differentiation (in vitro, single cell RNA-seq from Trapnell and Cacchiarelli et al., 2014 Nature Biotechnology)

Trapnell et al. analyzed in vitro human skeletal muscle myoblast (HSMM) differentiation. The HSMM dataset contained 372 single cells: 96 cells were from growth media (GM) before converting to differentiation media (DM). 96 cells were from 24 hours, 96 cells were from 48 hours and 84 cells were from 72 hours after converting to DM. Since Waterfall was developed to reconstruct the in vivo trajectory without any prior information annotated to cells, we merged all conditions, reconstructed the process and compared it with the Monocle predictions with or without the prior information.

Mapping and calculating normalized expression level

We used RSEM to calculate TPM values (Li and Dewey 2011).

Waterfall: Preprocessing

We performed unsupervised clustering and principal component analysis (Figure 5. 13A and C). We could determine that Branch 1 represented single cells from growth media (GM), where cells were highly proliferative with proliferation markers such as PCNA and CDK1 (Figure 5. 13B). On the contrary, Branch 2 represented single cells from a myogenic population in differentiation media (DM) with high myogenic markers such as ENO3 and MYOG (Figure 5. 13B). Since Branch 3 less likely represents myogenic population as opposed to the other two branches, we excluded Branch 3, thereby achieving a single myogenic trajectory from Branch 1 (GM; light purple, light green) to the Center (early stage in DM; red, yellow) and to Branch 2 (late stage in DM; dark orange, light orange) (Figure 5. 13C).

We found that, when given the temporal information, Monocle resulted in a similar trajectory, where Branch 1 largely corresponded to Group 2, Branch 2 largely corresponded to Group 1, and Branch 3 largely corresponded to Branch 3 (Data not shown). Trapnell et al. also excluded the Group 3 because they thought it represented contaminating fibroblasts.

Waterfall: Reconstructing trajectory and pseudotime

We reconstructed the MST by connecting k-means within the PCA plot (Figure 5. 14A) and assigned pseudotime to each individual single cell based on their relative location on the trajectory (Figure 5. 14B).

Waterfall: Gene expression analyses

Most genes showed expected expression patterns, suggesting the reconstructed trajectory reflected the myogenic differentiation process of HSMM (Figure 5. 15).

We sought to compare the performance of Waterfall without prior information, Monocle with prior information, and Monocle without prior information. First, we evaluated the gradual transition of the single cell transcriptomes along the trajectories as in Figure 5. 11 (Figure 5. 16A). Waterfall without prior information performed significantly better than Monocle with or without prior information. Second, we compared the predictions for marker genes from each method. Both

Waterfall and Monocle can predict novel myogenic markers by performing correlation analysis between gene expression and pseudotime progression. In order to test the performance of Waterfall and Monocle, we queried how well these algorithms could find the known myogenic genes (Figure 5. 16B).

3. Embryonic lung development (In vivo, single cell RNA-seq from Treutlein and Brownfield et al., 2014 Nature).

Mapping and calculating normalized expression levels

The dataset contained 199 single cell RNA-seq transcriptomes. We used RSEM to calculate TPM values (Li and Dewey 2011).

Waterfall: Preprocessing

We focused on the E18.5 single cells (79 cells). To select a major trajectory for the continuous process, we performed preprocessing steps of Waterfall.

We performed unsupervised clustering using Pearson correlation, followed by calculating distance matrix from the correlation matrix (Figure 5. 17A). For PC analysis, we followed the approach of Treutlein et al. (Figure 5. 17B) (Treutlein, Brownfield et al. 2014). First, we extracted 36 genes with the highest absolute loadings from each of the four highest contributing principal components. Second, we performed PCA only using these genes. According to Treutlein et al., E18.5 single cell dataset contained unrelated cell types including Clara cells and ciliated cells (Treutlein, Brownfield et al. 2014). Using *Scgb1a1*, *Krt15*, and *Foxj1*, we were able to determine the Clara cell- and ciliated cell-clusters (Figure 5. 17B-C). Furthermore, based on *Ager*, *Pdpn*, and *Aqp5* expression we determined the differentiating AT1 population. Based on *Sftpb* and *Sftpc* expression, we also identified the differentiating AT2 population (Figure 5. 17B-C). Consistent with Treutlein et al., the central population (BP) represented the stem cell population expressing both AT1 and AT2 markers at lower level than the both extremes (Figure 5. 17B-C).

Waterfall: Reconstructing trajectory and pseudotime

As opposed to Treutlein et al. where the authors used group-to-group comparisons to discover key developmental genes (Figure 5. 21), Waterfall allows for reconstruction of the predicted in vivo trajectory and analyses of gene expression at single cell resolution.

We eliminated the Clara cells and ciliated cells to reveal the two trajectories on the PCA plot (Figure 5. 18). We then performed Waterfall trajectory reconstruction for both AT1 and AT2 routes (Figure 5. 19A-B).

AT2 route

We reconstructed AT2 route using MST (Figure 5. 19A, top). We assigned pseudotime of each single cell based on its relative location on the trajectory (Figure 5. 19A, bottom).

AT1 route

We reconstructed the AT1 route using MST (Figure 5. 19B, top). We assigned pseudotime of each single cell based on its relative location on the trajectory (Figure 5. 19B, bottom).

Waterfall: Gene expression analyses

After reconstructing the trajectory, we could query the whole transcriptome profiles throughout the reconstructed in vivo AT1 and AT2 differentiation process (Figure 5. 20). It recapitulated the key discovery of Treutlein et al. with greater temporal resolution (Treutlein, Brownfield et al. 2014): The BP cells were located at the beginning of both pseudotime trajectories and cells in this cluster expressed low level of the both markers. BP cells gained AT1 markers (Figure 5. 20B, left) and lost AT2 markers along the AT1 route (Figure 5. 20A, left). In contrast, the BP cells gained AT2 markers (Figure 5. 20A, right) and lost AT1 markers along AT2 route (Figure 5. 20B, right).

Importantly, we observed that single cell level heterogeneity was resolved by the successful reconstruction of the trajectory. For example, we could appreciate the single cells expressing higher level of AT2 cell marker *Clic5* was placed earlier pseudotime points (asterisks in Figure 5. 20A) than the single cells expressing lower level of *Clic5* (arrows at the bottom of *Clic5* graph in Figure 5. 20A) along the AT1 trajectory, even though *Clic5* was not used for reconstructing the trajectory. These results suggest that the pseudotime reflects the whole transcriptomic shift during the developmental process.

In addition, we were able to appreciate the differential timing of gene expression transitions, which was revealed by the bimodal state prediction. Clustering-based differential expression analysis inevitably underestimates developmental heterogeneity (Figure 5. 21). Therefore, single cell resolution reconstruction and bimodal state prediction would be helpful in interpreting the single cell data in higher resolution and allow achieving additional insights.

4. B cell development (In vivo, single cell mass cytometry from Bendall et al. 2014 Cell)

Using their innovative single cell analysis software called Wanderlust, Bendall et al. successfully reconstructed the in vivo trajectory of in vivo B cell development from single cell mass cytometry data. Unfortunately, Wanderlust is not applicable for single cell RNA-seq analysis in its current form. Although Waterfall was developed for single cell RNA-seq, we found that it was also applicable for single cell mass cytometry dataset from Bendall et al.

Waterfall: Preprocessing

We herein present data analysis for the third biological replicate called “Sample C”, although we found that all four biological replicates resulted in consistent and reproducible conclusion. We used the 19 CD markers and HLA-DR, IgD, IgM, Kappa, and Lambda light chain genes as stratifiers for cellular states. Unsupervised clustering and PC analysis resulted in only a single trajectory (data not shown).

Waterfall: Reconstructing trajectory and pseudotime

Using the normalized marker gene expression table, we reconstructed the trajectory (Figure 5. 22, left) and then assigned pseudotime to each single cell data point based on the relative location on the reconstructed trajectory (Figure 5. 22, right).

Waterfall: Gene expression analyses

After reconstructing the trajectory, we performed expression analysis for all the markers in the dataset. Instead of performing HMM, we averaged expression levels within each bin and generated heat maps (Figure 5. 23A) as well as a polynomial regression fitting plot (Figure 5. 23B) to allow direct comparison with the Wanderlust output (Figure 5. 23C)(Bendall, Simonds et al. 2011). Most markers showed expected patterns of expression (Figure 5. 23A-C). Importantly, all four biological replicates generated comparable data without modification in parameters, suggesting the robustness of the Bendall et al. dataset as well as Waterfall (Data not shown).

5. Synthetic datasets

We generated two synthetic datasets. For synthetic dataset I, we allowed gene expression to change gradually over the sequence of single cell transition. For synthetic dataset II, we generated a more realistic condition that contains a branch coming out from a primary trajectory. We also introduced random stochasticity for dataset II.

Waterfall: Preprocessing

For synthetic dataset I, the variance contributions were distributed to four principal components (Figure 5. 24D). We demonstrated the gene expression level by the size of each single cell data point proportional to the gene expression level (Figure 5. 24D).

For synthetic dataset II, we introduced a branch as well as expression stochasticity (Figure 5. 24A, left and middle). Single cells originating from this branch were mixed with the single cells originating from the major trajectory of interest on the PCA plot (note the single cells with an asterisk at Figure 5. 25A, right). However, hierarchical clustering precisely determined the branch-derived single cells (Figure 5. 25B-C), so that we could exclude this branch from further analyses.

Waterfall: Reconstructing trajectory and pseudotime

For synthetic dataset I, the color-coded clusters lined up sequentially throughout the trajectory. K-mean clustering and MST successfully approximated the average trajectory (Figure 5. 26A), thereby generating a near linear graph without significant tangential offset (Figure 5. 26A).

For synthetic dataset II, after the branch was removed the major trajectory was successfully approximated by k-means and MST (Figure 5. 27A). When it was linearized and flattened as shown in the Figure 5. 27B, we were able to appreciate the stochasticity that we introduced forming the tangential distribution of the pseudotime progression (Figure 5. 27B).

Waterfall: Gene expression analyses

We checked a few gene expression profiles in the synthetic dataset I. As expected, the gene expression profile underwent smooth transition throughout the trajectory suggesting the reconstruction was precise (Figure 5. 28). Of note, the HMM was able to predict the transition points of each gene. Under naturally occurring conditions, these predictions could be validated by independent biological experiments.

Table 5. 1. List of datasets to which Waterfall was applied.

Dataset	Biological Process	Single Cell Modality	Bioinformatic Algorithms	Figures
Current study	<i>In vivo</i> adult neurogenesis	RNA-seq	Waterfall & Monocle	Figures 5. 1 –11
Trapnell et al., 2014	<i>In vitro</i> myogenesis	RNA-seq	Waterfall & Monocle	Figures 5 .12 – 15
Treutlein et al., 2014	<i>In vivo</i> lung development	RNA-seq	Waterfall	Figures 5. 16 – 20
Bendall et al., 2014	<i>In vivo</i> B cell development	Mass-cytometry	Waterfall	Figures 5. 21 – 22
Current study	Synthetic dataset I & II	–	Waterfall	Figures 5. 23 – 27

Table 5. 2 . Comparison between static and dynamic heterogeneity.

	Static/cell-type heterogeneity	Dynamic/temporal heterogeneity
Goals	Discovering novel cell types from a seemingly homogeneous population	Understanding the continuous biological process of one cell type
Dataset	Multiple cell types	One cell type at different temporal stages
Analytical approach	Clustering-based comparison (analogous to population RNA-seq analysis after FACS sorting)	Ordering single cell transcriptomes, based on the transcriptomic similarity (Waterfall or Monocle)
Examples	Unappreciated heterogeneity within sensory neurons in dorsal root ganglia (Usoskin, Furlan et al. 2015)	Molecular cascades throughout adult neurogenesis (Current study)

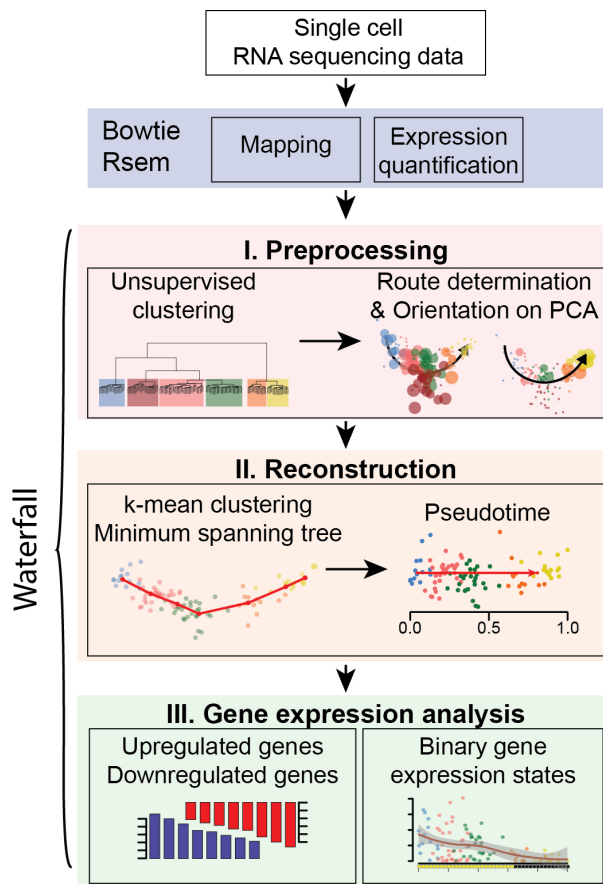


Figure 5. 1. Waterfall pipeline consists of three key steps of Pre-processing, Reconstruction, and Gene expression analysis.

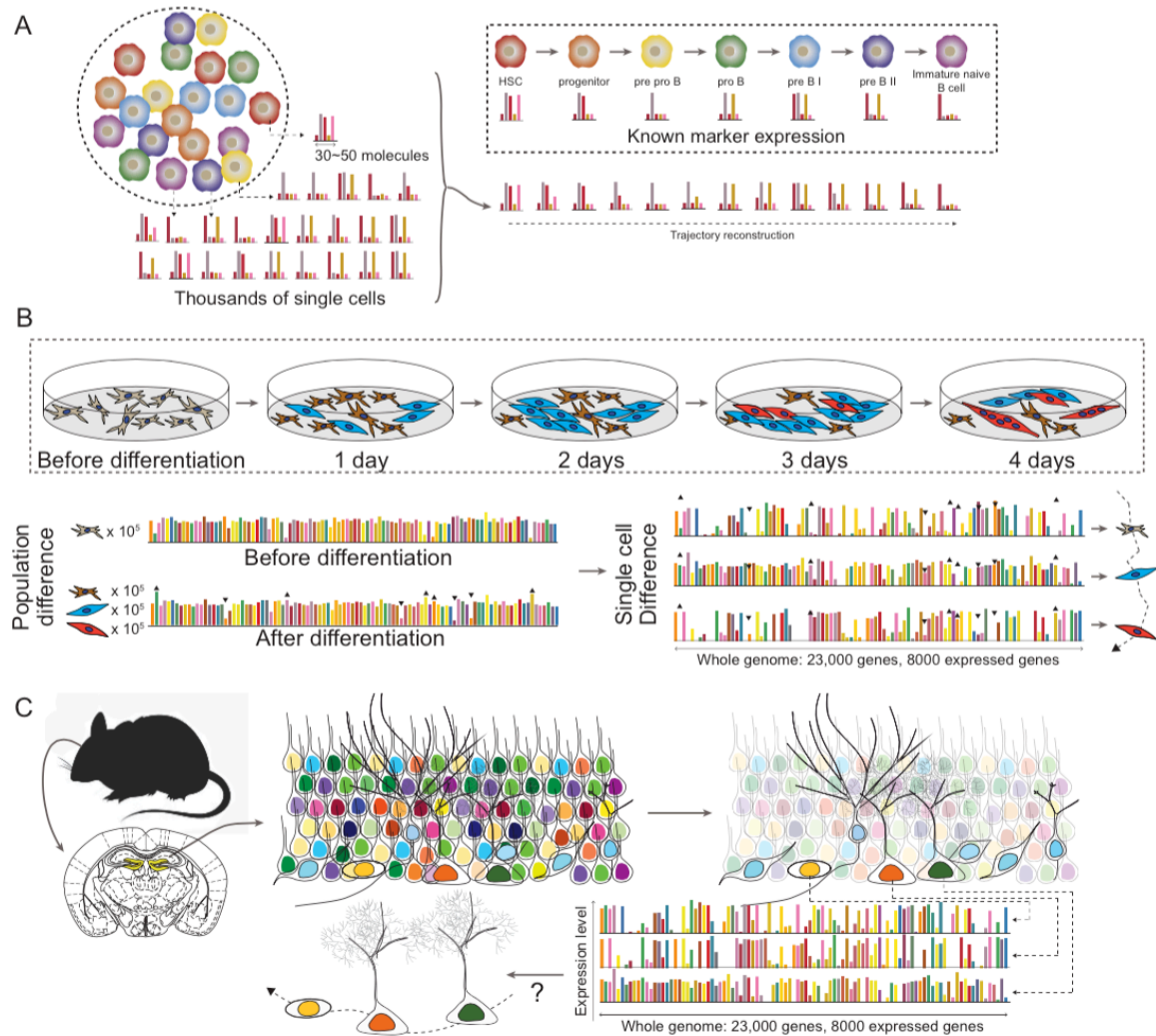


Figure 5. 2. Schematic comparison of the complexity between single cell qPCR or mass cytometry.

Single cell qPCR or mass cytometry generates data with 30~50 molecular markers for thousands of single cells (A). Using the known marker genes as stratifiers, we can order the single cells (A). For in vitro conditions, we can discover developmentally important marker genes by performing differential analysis using temporal information of single cells at population level (B). For in vivo conditions, however, we have a completely mixed population of differential developmental points (C). Thus, reconstructing the trajectory from the in vivo single cell is a significant challenge.

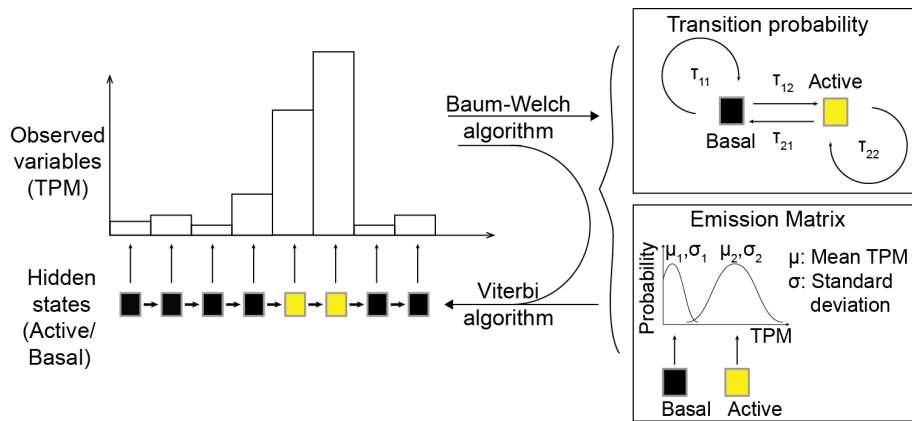


Figure 5. 3. Hidden Markov model (HMM). We predicted underlying states from gene expression (TPM) over pseudotime progression.

Baum-Welch algorithm predicts the most likely transition probability and emission matrix from observed variables (TPM). Viterbi algorithm uses observed variables (TPM) along with output from the Baum-Welch algorithm to predict hidden On/High and Off/Low states.

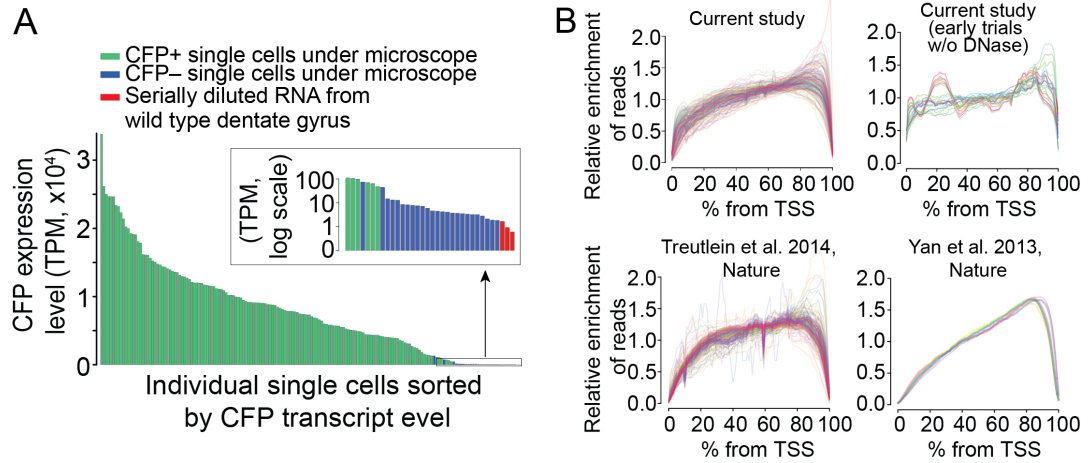


Figure 5. 4. Precision of CFPnuc⁺ single cell isolation (A) and minimal cDNA amplification related 3' bias (B).
Dnase treatment significantly augmented the quality of the sequencing (B).

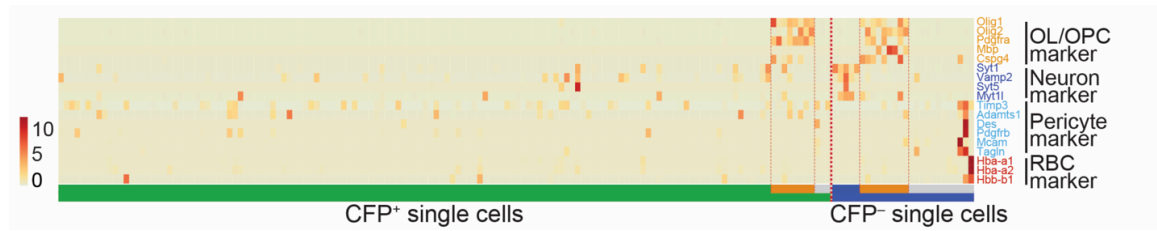


Figure 5. 5. Marker based cell type screening. Single cells expressing known markers for pericytes, oligodendrocyte (OL), oligodendrocyte progenitor cells (OPC) were excluded for the further analyses.

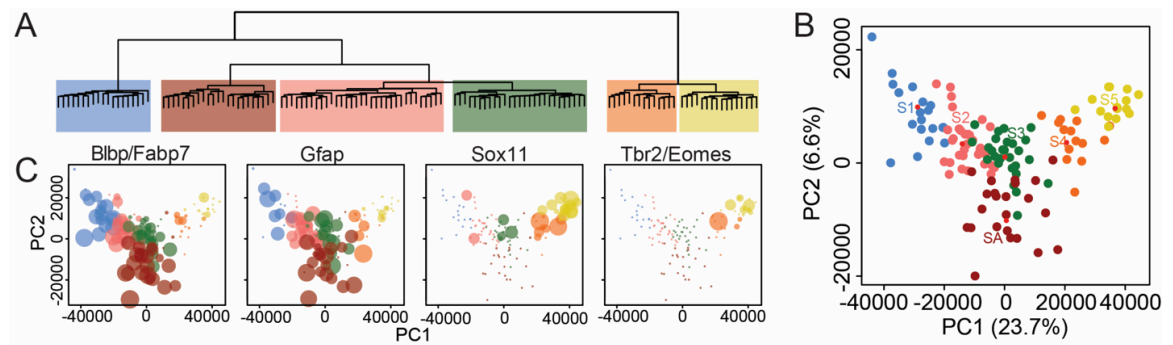


Figure 5.6. Preprocessing step of Waterfall - Unsupervised clustering (A), principal component analysis (B) and marker expression patterns (C) to help orient of data.

Color of each data point follows the same color code. Size of the each data point in C represents normalized expression level.

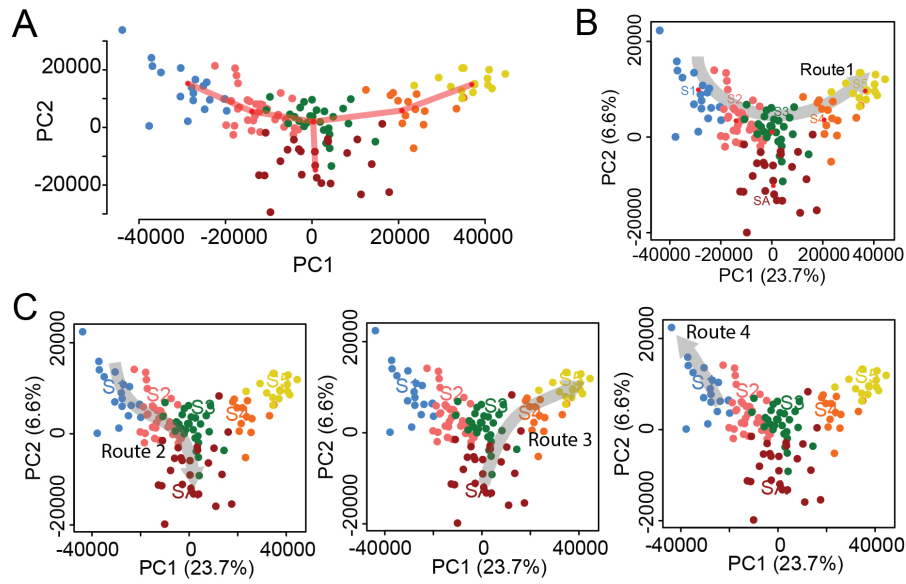


Figure 5. 7. Potential trajectories of adult neurogenesis dataset- route definition using minimum spanning tree with k-means clustering (A), the major neurogenic trajectory (B) and alternative routes to consider (C).

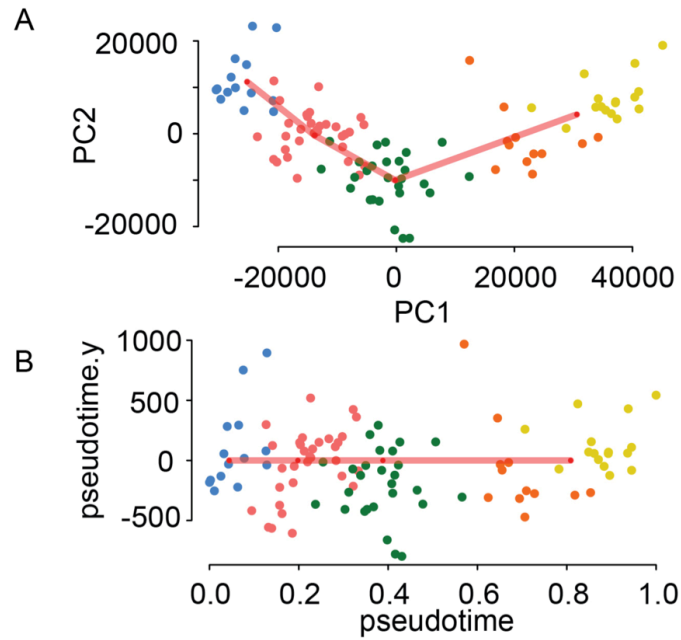


Figure 5. 8. Trajectory reconstruction step of Waterfall. The trajectory was built using MST connecting the five vertices from k-means clustering (A).

Pseudotime of each data points was assigned by its relative location when projected on to the reconstructed trajectory (B). The orientation was determined by Figure 5. 6C.

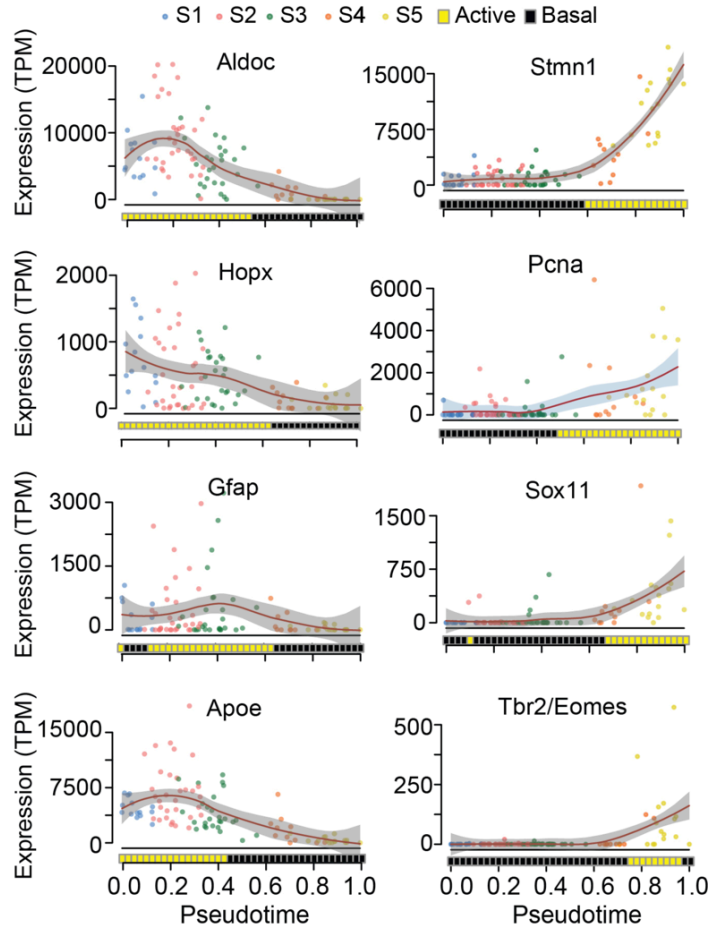


Figure 5. 9. Representative gene expression dynamics predicted by Waterfall. Shown are scatterplots of single cell gene expression levels.

We demonstrated local polynomial regression fitting plot with 95% confidence interval (red linear graphs and gray shades). HMM-predicted transcriptional states throughout pseudotime progression were demonstrated as block graphs at the bottom of each plot with on/high states as yellow blocks and off/low states as black blocks.

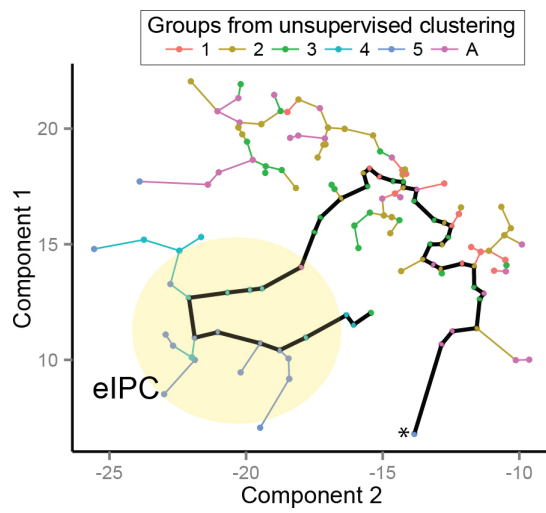


Figure 5. 10. Monocle-derived trajectory using top 1000 highly expressed genes.

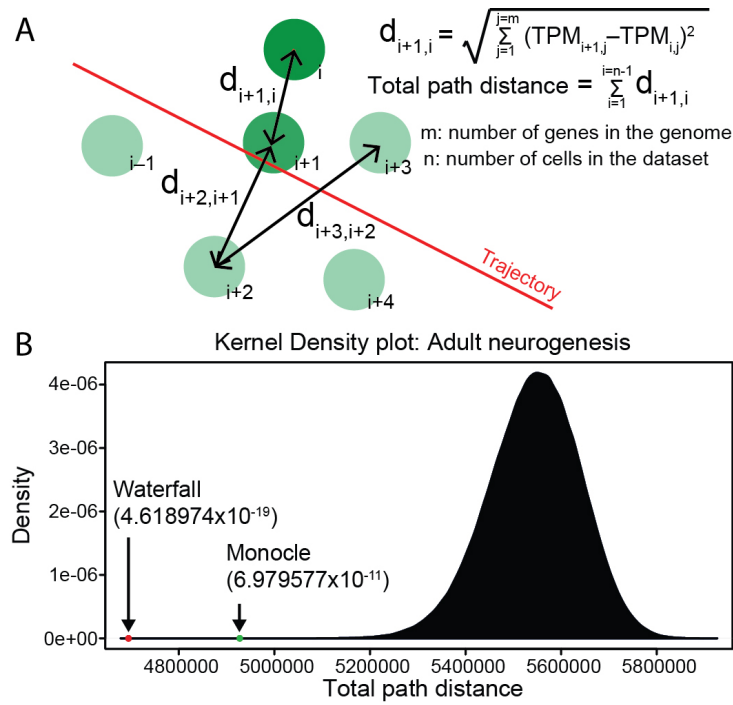


Figure 5. 11. The quantitative test for gradual transcriptomic transition of single cell trajectories.

(A) The distance between a pair of single cells is defined by the square root of the sum of square of each gene expression difference. Total path distance is defined by the sum of the distances between neighboring single cell pairs along the trajectory from the beginning to the end of the trajectory. (B) Kernel density plot represents the distribution of total path length from random ordering of single cells. Total path lengths of Waterfall- and Monocle-derived pseudotime order of single cells were marked by arrows. The cumulative probabilities for each predicted path were 4.6×10^{-19} for Waterfall and 7.0×10^{-11} for Monocle.

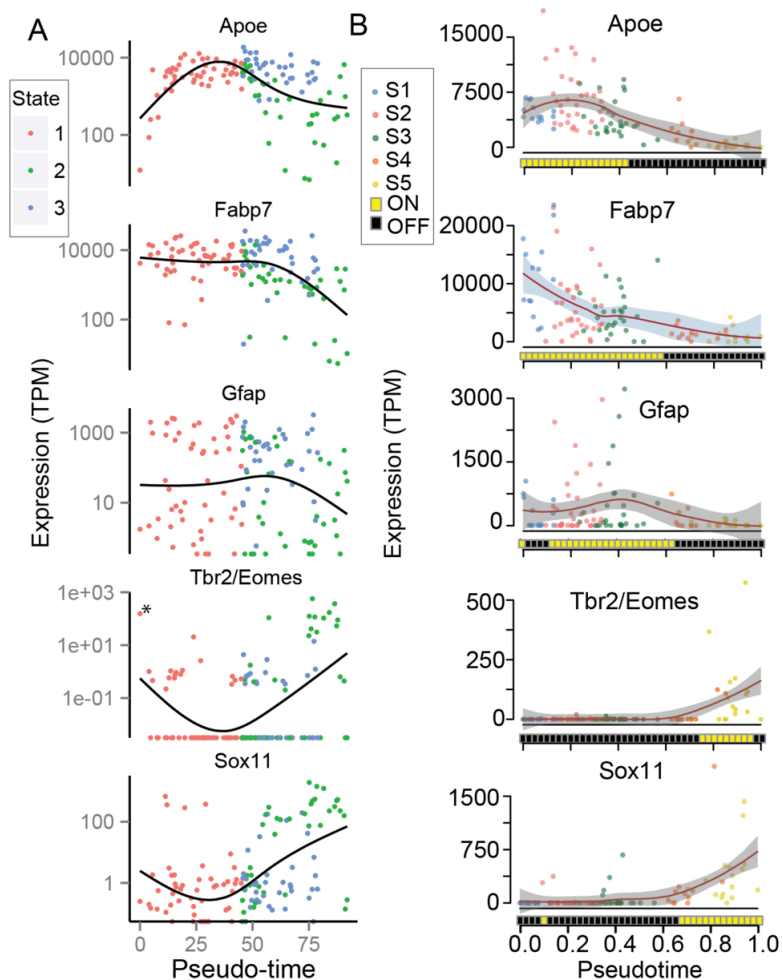


Figure 5. 12. Expression profiles of key neurogenic genes from Monocle (A) and Waterfall (B).

Note that Monocle represents the data as a logarithmic scale whereas Waterfall represents the data as a linear scale. The single cell with an asterisk at the Monocle predicted Tbr2/Eomes profile indicates one of misplaced cells, which was also marked at Figure 5. 10.

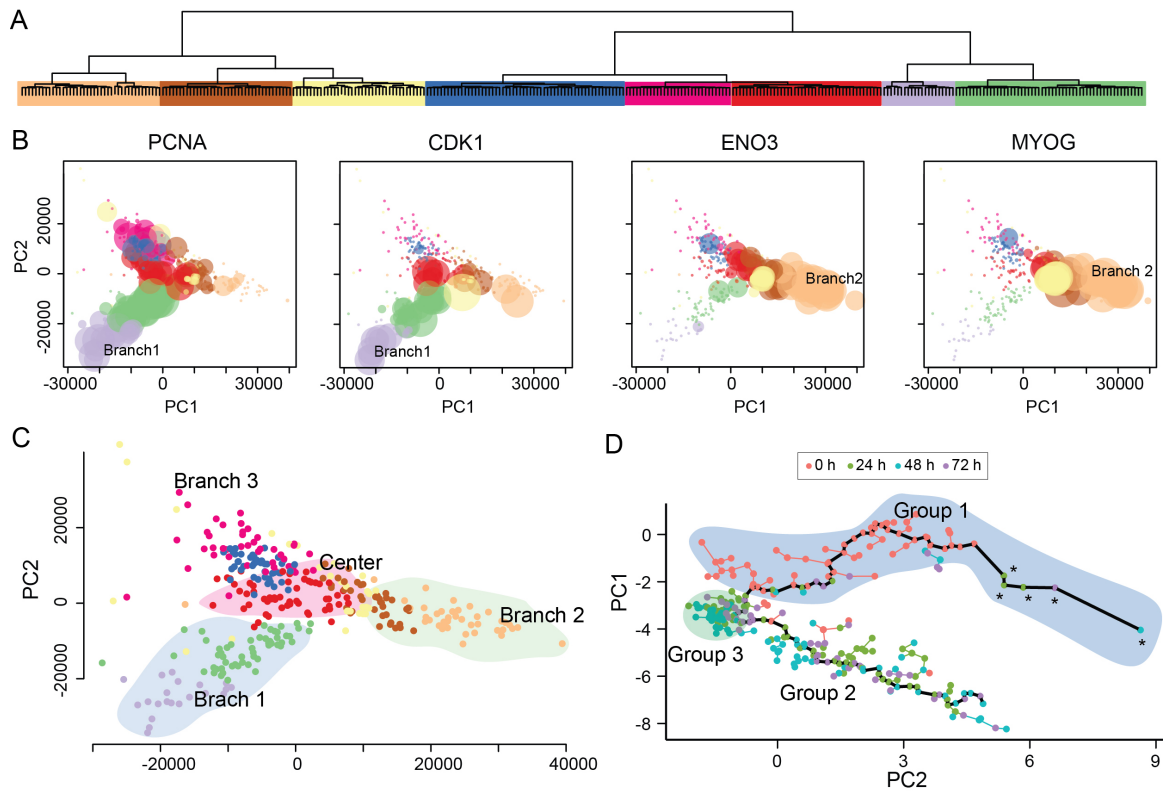


Figure 5.13. Preprocessing steps of Waterfall: Unsupervised clustering (A), marker gene expression profiles (B), Waterfall trajectory determination without prior temporal information (C); and Monocle trajectory with prior temporal information as a comparison (D).

The color in B and C follows the color code appearing in A. The size of each data point at B represents normalized expression value.

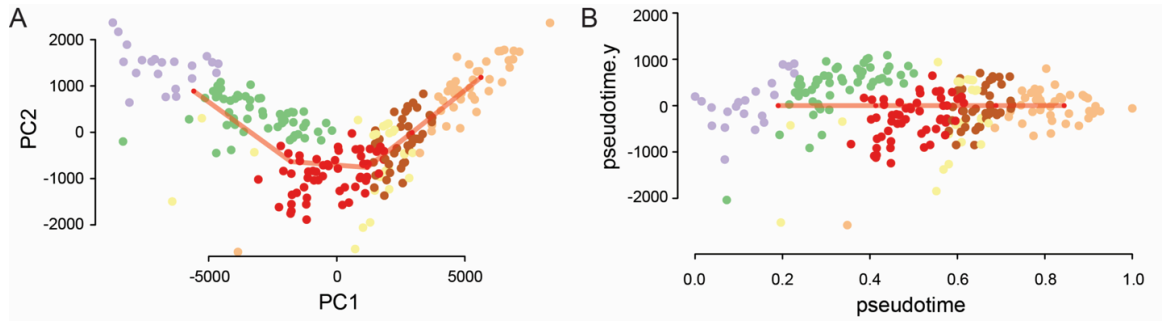


Figure 5. 14. Pseudotime reconstruction of selected route of in vitro HSMM differentiation. MST connected k-means to generate trajectory (A).

Pseudotime was determined based on the relative location of single cells when projected onto the reconstructed trajectory (B). The color of each data point follows the color code at Figure 5. 13 A.

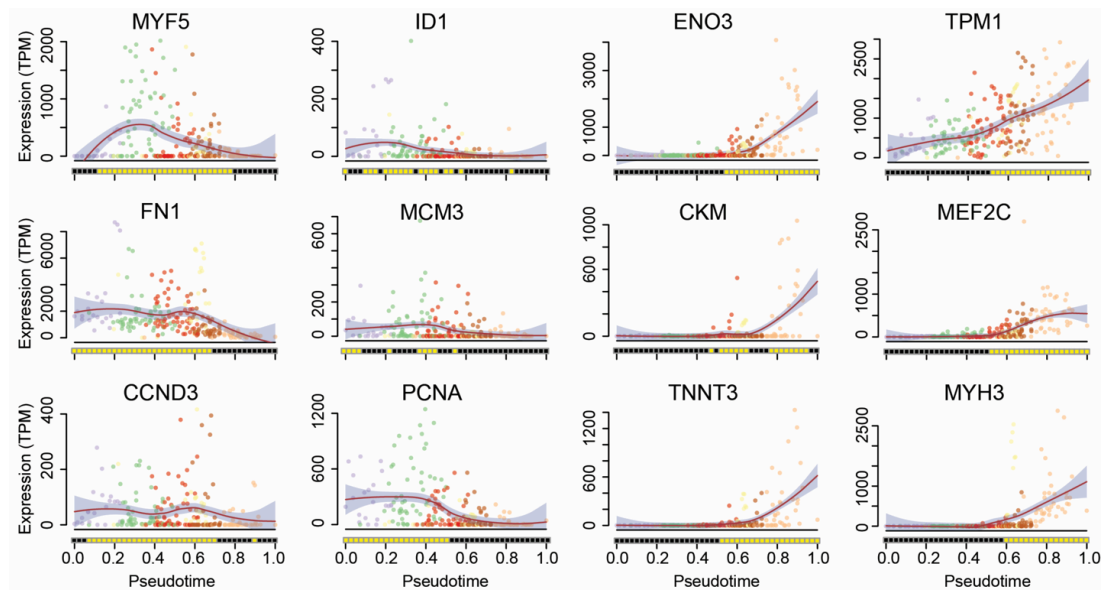


Figure 5.15. Gene expression dynamics over the reconstructed HSMM differentiation process.

Shown are the local polynomial regression fitting plots with 95% confidence interval (red linear graphs and gray shades), and HMM-predicted transcriptional states (block graphs at the bottom; on/high-yellow, off/low-black). Proliferative marker genes or primitive muscle developmental genes were highly expressed at earlier pseudotime points (left two columns) whereas skeletal muscle markers were highly expressed at later pseudotime points (right two columns).

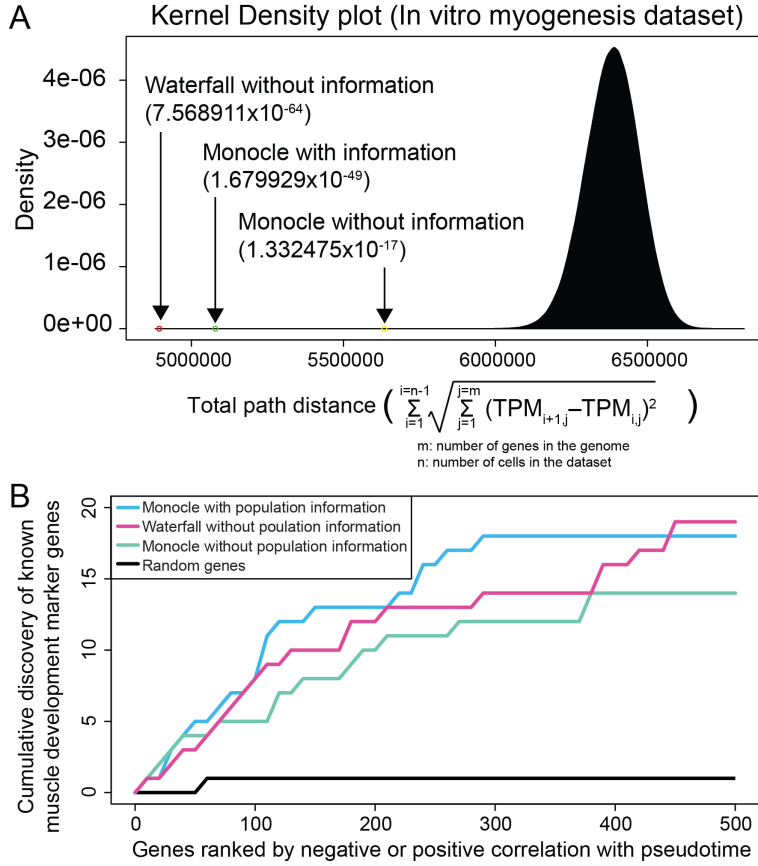


Figure 5. 16. Comparison between Waterfall and Monocle for in vitro myogenesis dataset.

(A) Kernel density plot represents the distribution of the total path length of random order of single cells. Arrows indicate the total path length of pseudotime order of single cells from Waterfall without prior information, Monocle with prior information and Monocle without prior information. The cumulative probabilities for each predicted path were 7.6×10^{-64} for Waterfall without prior information, 1.7×10^{-49} for Monocle with prior information and 1.3×10^{-17} for Monocle without prior information. (B) Discovery of myogenic genes using different methods of trajectory building. We counted the number of known myogenic genes within the top 500 genes with highest positive or negative absolute correlation with pseudotime progression predicted from (1) Monocle with prior information (light blue), (2) Waterfall without prior information (pink), (3) Monocle without prior information (light green) and (4) randomly picked genes (dark gray). List of markers for muscle development was from Trapnell et al. (Trapnell, Cacchiarelli et al. 2014) and http://www.rndsystems.com/molecule_group.aspx?g=821.

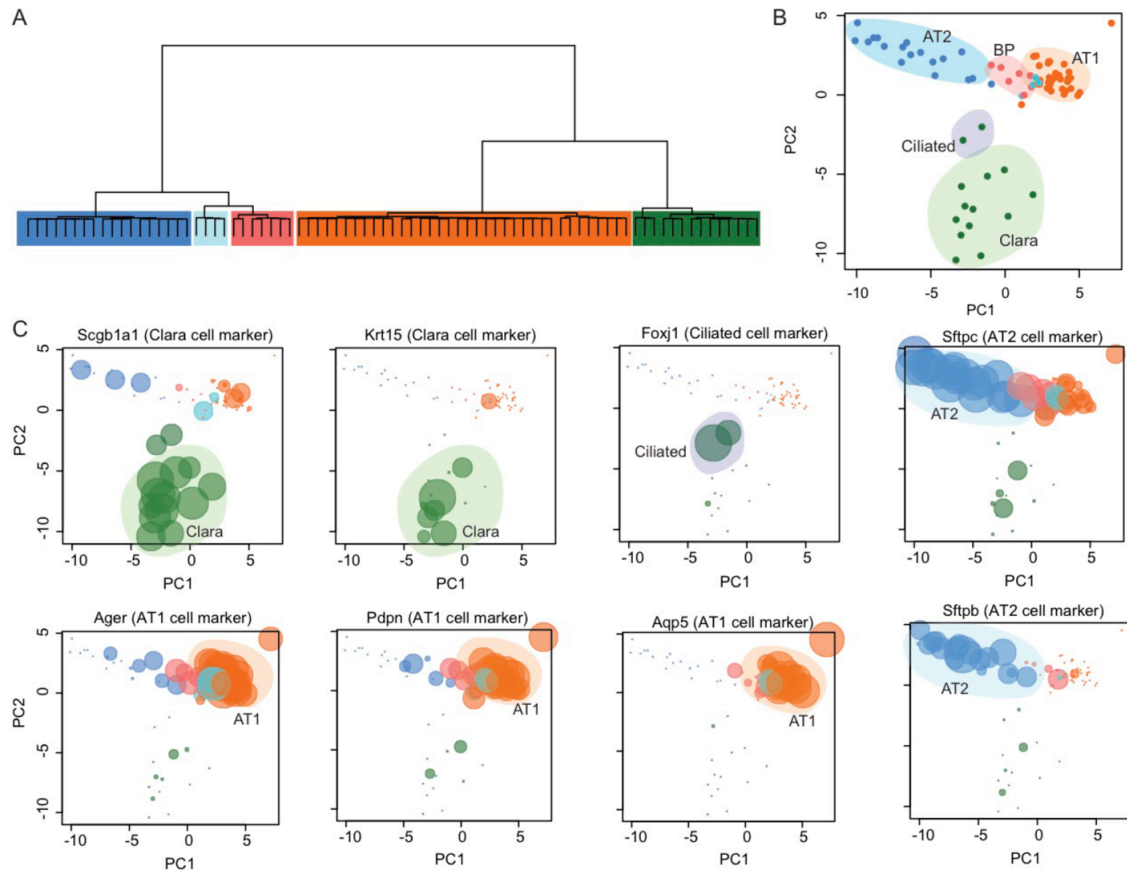


Figure 5.17. Preprocessing step of Waterfall for embryonic lung development: including unbiased clustering (A), PCA plot using selected genes from the four major PC axes (B), and marker gene expression profiles to orient and define a single route of interest (C).

Size of each data point in C represents normalized gene expression and the color follows the color codes in A.

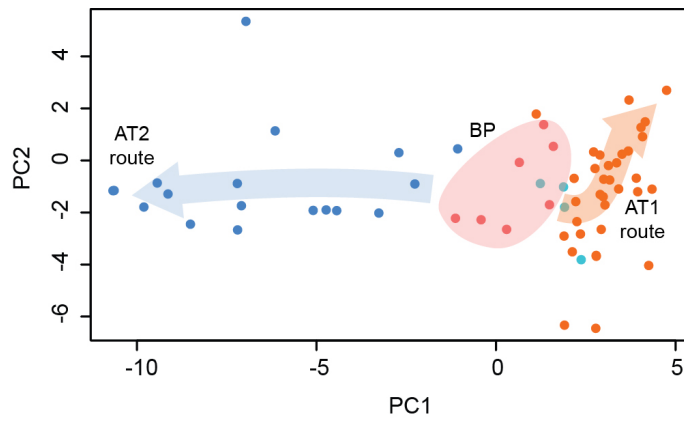


Figure 5. 18. PCA plot after removing Clara/ciliated cell group.

AT1 route and AT2 route are shown with orange and blue arrows, respectively.

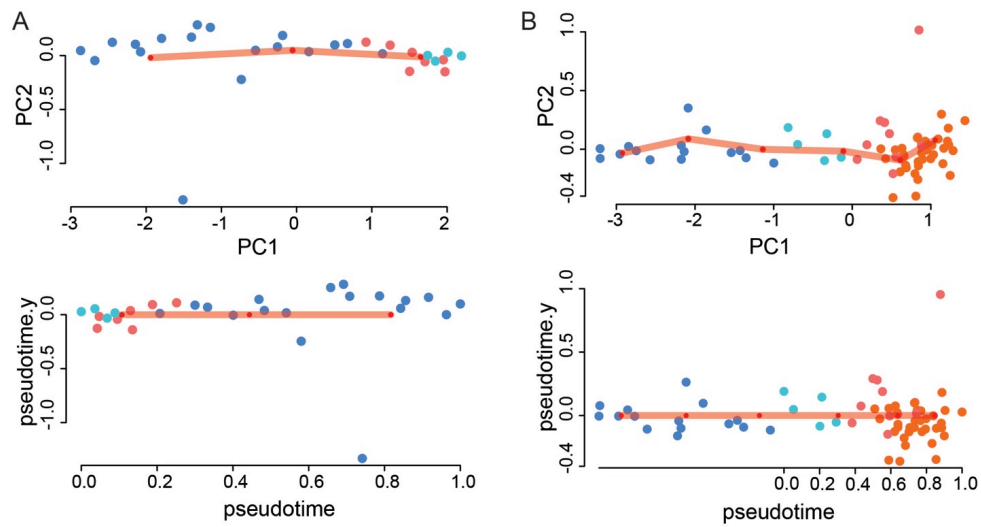


Figure 5.19. Reconstructing continuous trajectory and assigning pseudotime for AT2 (A) and AT1 differentiation process (B).

The vertices from the k-means clustering were connected by MST (A, top; B, top). Pseudotime was determined by the relative location of single cells when projected on to the trajectory (A, bottom; B, bottom). The color of each data point follows the color code appearing in Figure 2. 3

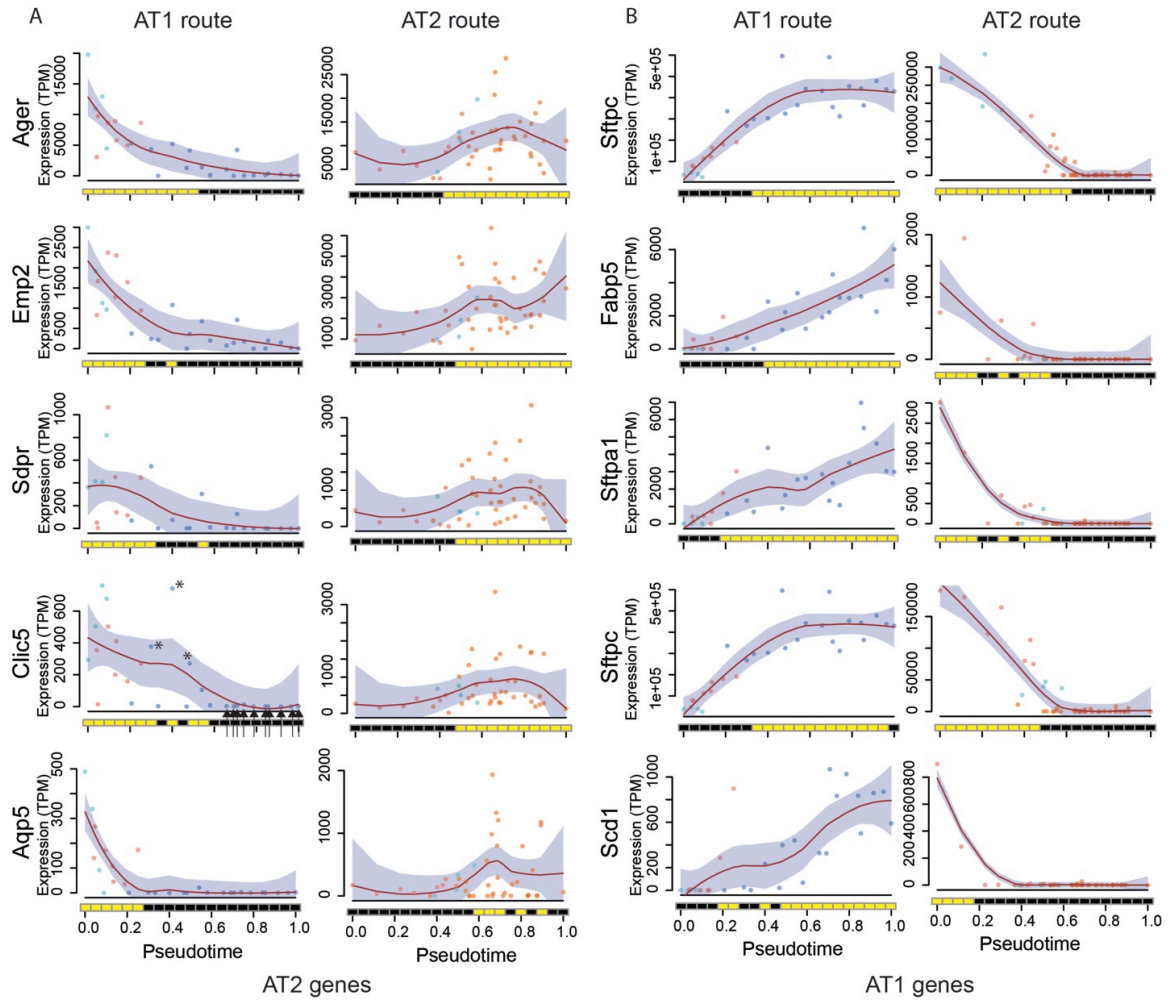


Figure 5.20. Waterfall gene expression prediction profile for AT2 markers (A) and AT1 markers (B).

We showed the local polynomial regression fitting plot with 95% confidence interval (red linear graphs and gray shades), and HMM-predicted transcriptional states (block graphs at the bottom; on/high-yellow, off/low-black). BP cells gained AT1 markers (B, left) and lost AT2 markers through AT1 route (A, left). In contrast, BP cells gained AT2 markers (A, right) and lost AT1 markers through AT2 route (B, right). Even though single cells indicated by the arrows and single cells indicated by asterisks at Clic5 expression profile were classified as a same group by unsupervised clustering, Waterfall successfully placed them accordingly, further highlighting the greater resolution of Waterfall. Of note, Clic5 was not one of the genes that was used for trajectory reconstruction.

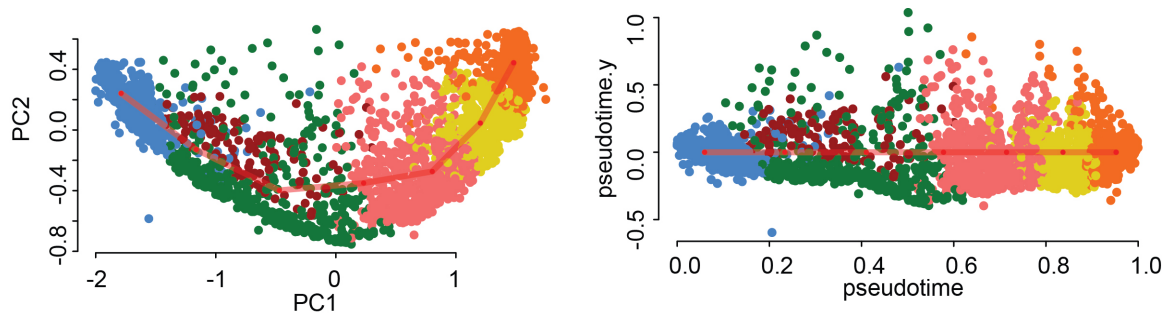


Figure 5.22. Waterfall trajectory reconstruction for the “Sample C” of Mass cytometry B cell development dataset. Vertices from the k-means clustering were connected by MST (left).

Pseudotime of each data points was assigned by the relative location when projected on to the reconstructed trajectory (right). The color of each data point follows the color codes from the unsupervised clustering (data not shown).

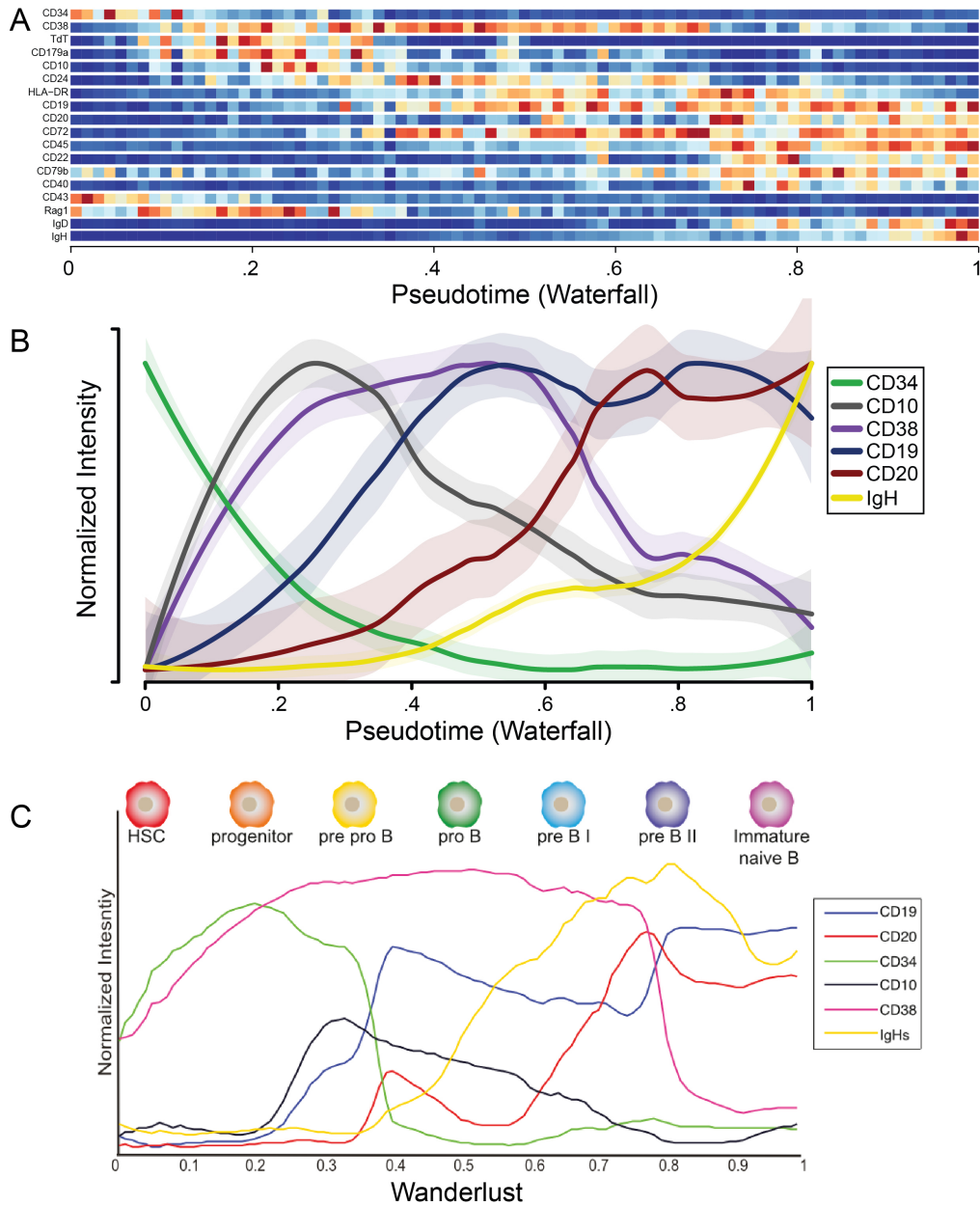


Figure 5.23. Waterfall predicted gene expression profiles along the B cell developmental trajectory from the “Sample B” (A and B) and the Wanderlust predicted profile (C).

(A) The reconstructed trajectory by Waterfall was divided into 80 bins and the mass cytometry-measured antigen abundance was averaged within each bin. (B) We selected the key marker antigens for B cell development from Waterfall analysis. We showed the local polynomial regression fitting plot with 95% confidence interval (shaded region). (C) Key marker antigen profiles for B cell development, taken from Bendall et al. (Bendall, Davis et al. 2014).

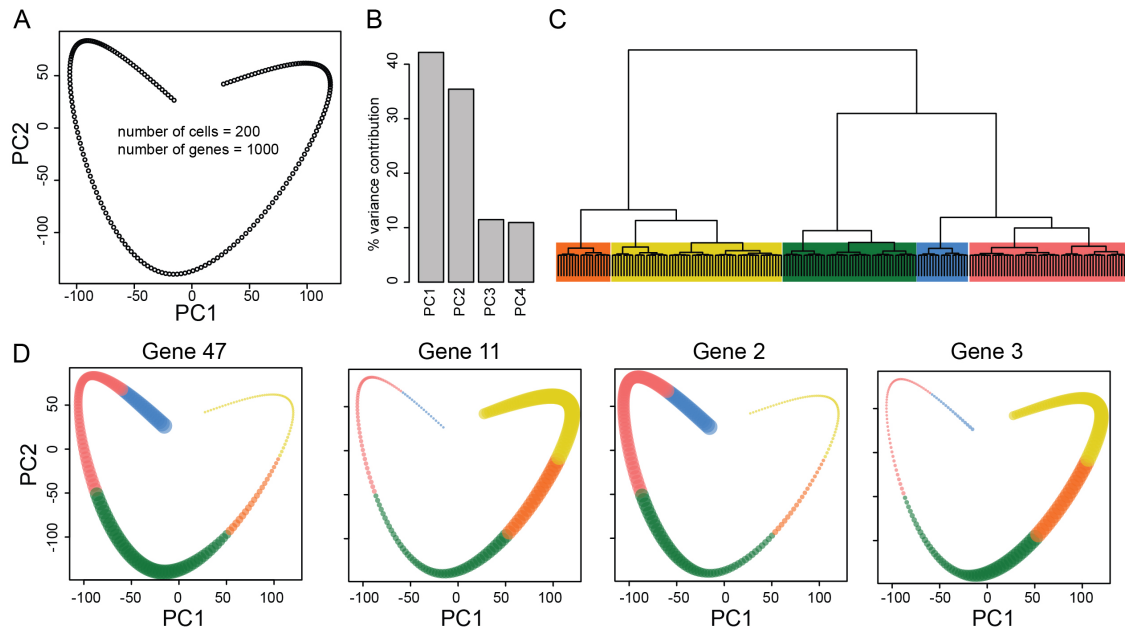


Figure 5. 24. Preprocessing of Waterfall analysis for synthetic dataset I.

The dataset consisted 1000 hypothetical genes gradually changing throughout preset order of 200 single cells, thereby generating a linear trajectory on the PCA plot (A). Four principal components could explain the variance (B). Although unsupervised clustering resulted in clusters, the clusters were not robust due to the continuity nature of the dataset (C). Examples of hypothetical gene expression were shown to mimic the preprocessing step of Waterfall (D). The size of the each data point is proportional to the gene expression level and the color of the each data point follows the color code in C.

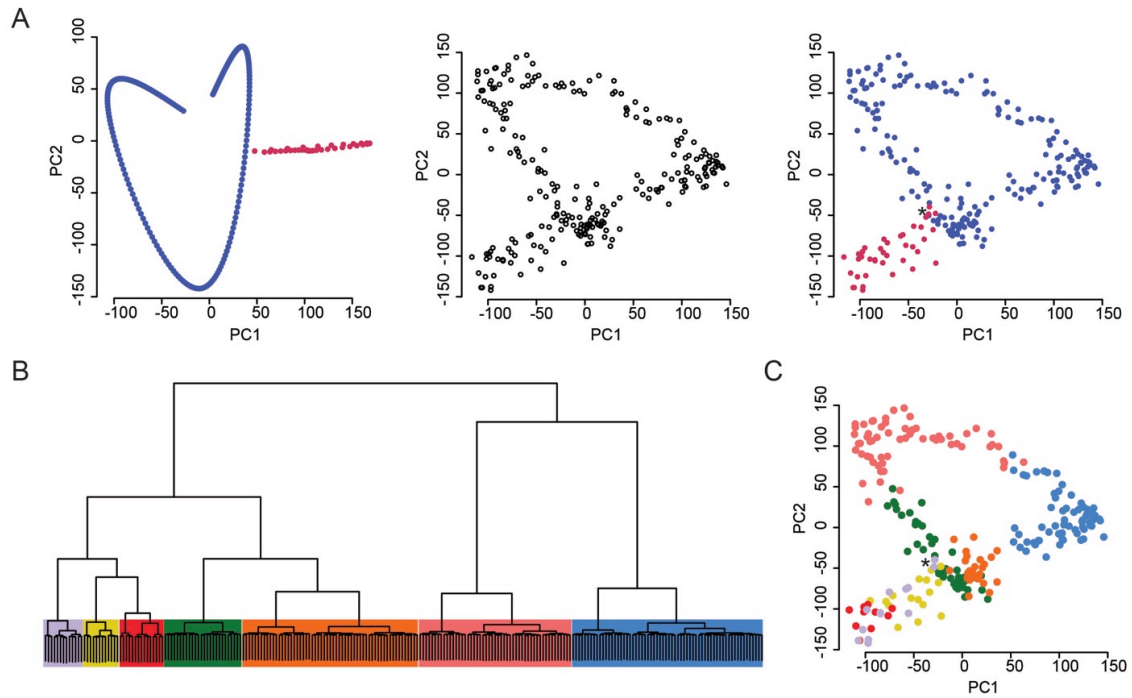


Figure 5. 25. Preprocessing of Waterfall analysis for the synthetic dataset II.

The linear dataset with a branch (A, left) was introduced with stochasticity (A, middle). It was ambiguous where the single cells close to the branching point belong to (A, right, asterisk), but unsupervised clustering (B) successfully determined them (C, asterisk). We cut out single cells that belong to the branch for the further analysis, which is a simulation for preprocessing for the real dataset appearing at Figure 5. 7 and Figure 5. 13.

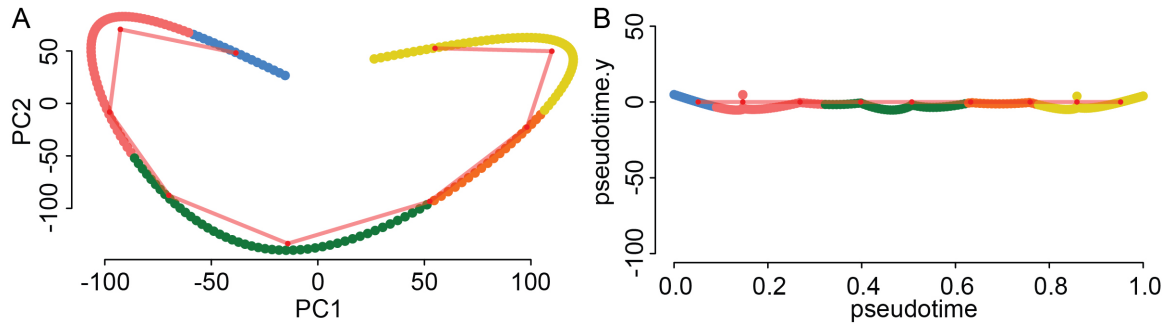


Figure 5. 26. Waterfall reconstructed trajectory and pseudotime assignment for synthetic dataset I.

Vertices from the k-means clustering were connected by MST (A). Pseudotime of each data point was assigned by the relative location when projected on to the reconstructed trajectory (B). The color of each data point follows the color codes appearing in Figure 5. 24C.

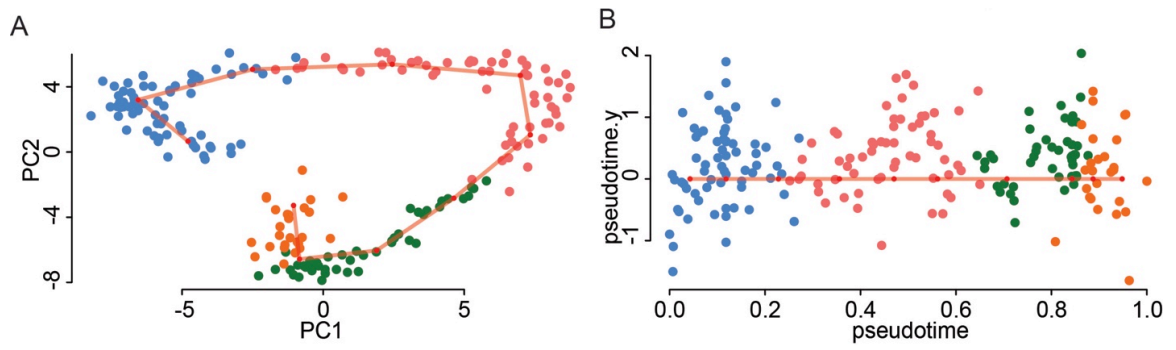


Figure 5.27. Waterfall reconstructed trajectory and pseudotime assignment for synthetic dataset II.

Vertices from the k-means clustering were connected by MST (A). Pseudotime of each data points was assigned by the relative location when projected on to the reconstructed trajectory (B). The color of each data point follows the color code in Figure 5.25B.

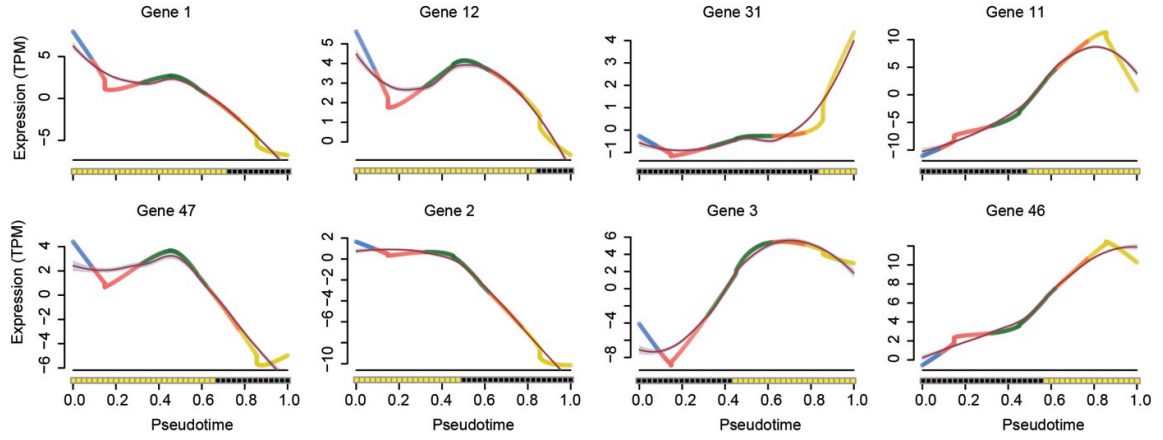


Figure 5.28. Gene expression profiles for synthetic dataset I, predicted by Waterfall.

The color of each data point follows the color code in Figure 5.24C. We showed the local polynomial regression fitting plot (red linear graphs), and HMM-predicted transcriptional states (block graphs at the bottom; on/high-yellow, off/low-black).

Chapter 6. Experimental procedures⁶

⁶ This chapter is based on Shin, J., et al. (2015). "Single-cell RNA-seq with waterfall reveals molecular cascades underlying adult neurogenesis." Cell Stem Cell **17**(3): 360-372.

Preparation of Individual Cells from Adult Mouse Dentate Gyrus

Homozygous transgenic mice expressing nuclear localized CFP (CFP^{nuc}) driven by the Nestin promoter (Encinas et al., 2006) were used for all single-cell RNA-seq experiments. Mice were euthanized by cervical dislocation, and brains were immediately immersed into ice cold Dulbecco's Phosphate-Buffered Saline (DPBS, Corning). All procedures were performed with approved protocols in accordance with institutional animal guidelines.

The dissected dentate gyrus was incubated in Hibernate A (BrainBits) containing papain (100 U; Sigma) and RNase-free DNase I (100 units; NEB) at 37°C for 18 minutes with intermittent flicking. The tissue was triturated into individual cell suspension by 1 ml pipette (Denville Scientific). Enzymes and cellular debris were removed with multiple rounds (~4-5 times) of mild centrifugation at 200g and washing with Hibernate A minus Ca²⁺ and Mg²⁺ (BrainBits). The individual cell suspension was plated onto a glass bottom plate (MatTek) and picked up using glass pipettes (World Precision Instruments) under a fluorescent microscope. The glass tip was broken into the bottom of each PCR tube containing water (2.4 μ l) with RNase-free DNase I (0.2 μ l; NEB) and Murine origin RNase inhibitor (0.25 μ l; NEB). Importantly, the addition of DNase I significantly improved the quality of the data by removing contamination from the random amplification of genomic DNA (See Chapter 5. Single cell analysis with Waterfall and Figure 5. 4B).

Library Preparation and Sequencing

cDNA amplification followed the previously published SMART protocol (Ramskold et al., 2012). Briefly, the DNase I was first inactivated by increasing the temperature (75°C for 10 minutes) and samples were then stored on ice. Custom designed 2A oligo 1 μ l (12 μ M, Integrated DNA Technologies, sequence shown in Figure 2. 2A) was added and annealed to the polyadenylated RNA by increasing temperature (75°C for 3 minutes) and quenching on ice. A mixture of 2 μ l Superscript II First-Strand Buffer (5X, Invitrogen), 1 μ l custom designed TS oligo (12 μ M, Integrated DNA Technologies, Figure 2. 2A), 0.3 μ l MgCl₂ (200 mM, Sigma), 0.5 μ l RNase inhibitor (Neb), 1 μ l dNTP (10 mM each, Thermo), 0.25 μ l DTT (100 mM, Invitrogen), and 1 μ l Superscript II (200 U/ μ l, Invitrogen), were incubated at 42°C for 90 minutes, followed by enzyme inactivation at 75°C for 10

minutes. A mixture of 29 μ l Water, 5 μ l Advantage2 taq polymerase buffer, 2 μ l dNTP (10 mM each, Thermo), 2 μ l custom designed PCR primer (12 μ M, Integrated DNA Technologies, Figure 2. 2A), 2 μ l Advantage2 taq polymerase was directly added to the reverse transcription product and the amplification was performed for 19 cycles. The amplification product was purified using Ampure XP beads (Beckman-Coulter). Library preparation was performed using Ovation Ultralow library systems (Nugen inc). Libraries were multiplexed and sequenced using Illumina Hiseq 2500 (Illumina Inc) (Table 2. 1).

Bioinformatic Analyses

Mapping and calculating gene expression levels: Raw reads were first trimmed for their Illumina adapter sequences using Trimmomatic (Bolger et al., 2014). Custom R codes were used to further trim the 5' TS oligo sequences, and the 3' primer sequences. RSEM (Li and Dewey, 2011) was used to map and calculate gene expression levels represented as transcripts per million (TPM). The reference genome was modified to include chrC, which contained the sequence of part of Nestin enhancer followed by eCFP transcripts, and reconstructed from sequencing reads. We used following parameters: `rsem-calculate-expression -p 12 --fragment-length-mean 500 $input.fastq $rsem_ref $cell_id`. For most of the downstream analyses, we used a table with single cells at the columns and the genes at the rows. For the purpose of visualization (Figure 2. 1C), reads were separately processed using bowtie (Langmead and Salzberg, 2012) and tophat (Trapnell et al., 2009) with the following option: `tophat2 -p 8 -N 2 -o $cell_id -g 1 -G $gtf --transcriptome-index $transcriptomeindex $bowtie2index $input.fastq`.

Waterfall 1--pre-processing: Waterfall input is an expression matrix from RSEM after eliminating outliers (Figure 2. 2C). Unsupervised clustering was performed using a distance matrix based on Pearson correlation between each pair of single cells (Figure 2. 3A). We defined the neurogenic trajectory on the PCA plot and determined the direction using known markers such as Sox11, Tbr2, Blbp, and Gfap. Expression levels represented on the PCA plot at Figure 2. 4B were based on custom normalization using sine normalization.

Waterfall 2 - Building an in vivo trajectory: We used custom R codes to determine pseudotime for each single cell on the trajectory (See Chapter 5. Single cell analysis with Waterfall). Briefly, we performed parametric PCA, and extracted k-means from the distribution of single-cell transcriptomes (Figure 2. 6A). We generate an unbiased trajectory by connecting k-means centers using a minimum spanning tree (MST) algorithm (Paradis et al., 2004). First, we set zero for the origin of the continuous trajectory, determined by pre-processing. Second, we assigned locations to individual cell data points on the trajectory. We assigned each cell to the closest MST segment (lines between k means) or vertex (k-mean) with a single perpendicular projection. Third, we straightened all the segments (Figure 2. 6A.d) into one horizontal line, and determined the relative order of the assigned locations of single cell data points on the trajectory. Pseudotime values ranged from 0 (at the origin) to 1 (at the end) (Figure 2. 6.d).

Waterfall 3 - Gene expression analysis by Hidden Markov model: We used custom R codes to apply a Hidden Markov model (HMM) to predict gene expression states throughout pseudotime. Briefly, we divided pseudotime into 40 bins, each of which contained an average of 2.5 single cells. We averaged the expression level within each bin and assigned the expression values to observed variables for HMM. We used Baum-Welch algorithm to extract the most probable emission probabilities and transition probabilities. Using the output from Baum-Welch algorithm along with observed variables, we applied the Viterbi algorithm to predict binary gene expression states (Figure 2. 6B).

Functional gene expression analysis: We calculated the Spearman correlation coefficient between pseudotime points and each gene's expression TPM values. Genes with relatively high Spearman correlations were defined as UP genes and genes with relatively low correlations were DOWN genes and the highest and lowest 1,000 genes defined as UP¹⁰⁰⁰ and DOWN¹⁰⁰⁰, respectively (Table 3. 1). A small subset of the UP¹⁰⁰⁰ and DOWN¹⁰⁰⁰ genes with low average expression values (< 50 TPM) and low coefficient of variation (< 1.95) were from repeat elements within their exons and excluded from downstream analyses. Raw mapping profiles of all genes shown in pseudotime figures were closely inspected to rule out false positives. We identified transcription factors using public databases (Zhang et al., 2012). We used GO (Ashburner et al., 2000), Kyoto Encyclopedia of Genes and Genomes (KEGG) (Kanehisa and Goto, 2000) and

wikipathway (Pico et al., 2008) for functional gene ontology analyses, along with the R bioconductor package (Gentleman et al., 2004) or Cytoscape software (Shannon et al., 2003).

For alternative functional gene expression analysis (Figure 3. 8), we first divided the entire transcriptome dataset into three equal groups based on their Spearman correlation to pseudotime: positively correlated, uncorrelated, and negatively correlated. We then evaluated the proportion of positively correlated genes versus negatively correlated genes within each functional entity from an independent functional annotation database wikipathway rather than the KEGG pathway database. If a functional entity contained a disproportionally larger number of up-regulated genes than down-regulated genes, we considered the functional entity to be generally activated, and, conversely, a disproportionally larger number of down-regulated genes indicated that the pathway was generally inactivated over time.

Validation with in situ Database Comparison, Immunohistology, Genetic Labeling and Electrophysiology

We validated our NPC-enriched genes by the Allen mouse brain atlas in situ hybridization dataset (Lein et al., 2007) (Table 2. 3). We inspected the gene expression patterns within the adult dentate gyrus at sagittal views. Genes with clear and relatively even distribution within the SGZ were determined to be “SGZ enriched”, whereas genes with subtle or scattered enrichment within the SGZ were determined to be “ambiguous”. Genes without any enrichment at the SGZ were determined to be “not SGZ enriched”.

Adult *Nestin-GFP^{cyto}* animals (Encinas et al., 2011) were used for immunohistochemical validation. *Hopx-CreER^{T2 f/+};;mT/mG^{f/+}* mice were generated by crossing *Hopx-CreER^{T2 f/+}* (Takeda et al., 2011) (Strain: *Hopx^{tm2.1(cre/ERT2)Joe}/J*, Jackson Labs Stock: 017606) with the *mT/mG^{f/+}* reporter line (Strain: B6.129(Cg)-Gt(ROSA)26Sor^{tm4(ACTB-tdTomato,-EGFP)Lo}/J; Jackson Labs Stock: 007676). Tamoxifen (62 mg/ml; Sigma; T5648) was prepared in a 5:1 ratio of corn oil/ethanol and heated to 37°C and mixed. Eight week-old *HopX-CreER^{T2 f/+};;mT/mG^{f/+}* animals were injected intraperitoneally with 124 mg/kg tamoxifen. Animals were perfused and extracted brains were placed in 4% paraformaldehyde overnight at 4°C, and then 30% sucrose for 48 hours at 4°C before coronal sections (45 µm) were cut. Immunohistology was performed using antibodies as previously described (Bonaguidi et al., 2011). The following antibodies were used: Aldoc (1:200, goat;

Cat#SC12065; Santa Cruz), GFAP (1:2000, rabbit; Cat#Z0334; DAKO), GFP (1:1000, chicken; Cat#GFP-1020; Aves), GFP (1:1000, goat; Cat#600-101-215; Rockland), Nestin (1:500, chicken; Cat#NES; Aves), PCNA (1:2000, rabbit; Cat#ab18197; Abcam), PCNA (1:500, goat; Cat#SC9857; Santa Cruz), Stmn1 (1:200, rabbit; Cat#ab24445; Abcam), Tbr2 (1:1000, rabbit Cat#Ab23345; Abcam). GFP cells were identified with an Axiovert 200M microscope (Zeiss) and then acquired as z-stacks on Zeiss 710 single-photon confocal microscope using 40X or 63X objectives. For quantification of Stmn1, Aldoc and PCNA expression in *Nestin-GFP^{cyto}* mice, Z-stacks were acquired from 3 animals. Images were analyzed using Imaris 7.1.1 (Bitplane). RGLs were identified by their radial process and soma situated in the SGZ and IPCs were identified by their small soma and tangential process as previously described (Bonaguidi et al., 2011).

Adult *nestin-GFP^{cyto}* transgenic mice were used to validate expression of functional glutamate receptors on NSCs. GFP⁺ radial glia like NSCs in slices prepared acutely from adult animals were recorded by whole-cell patch-clamp upon puffing of AMPA or NMDA in the presence or absence of antagonists as previously described (Song et al., 2012).

References

Arnoldini, M., et al. (2014). "Bistable expression of virulence genes in salmonella leads to the formation of an antibiotic-tolerant subpopulation." PLoS Biol **12**(8): e1001928.

Bendall, S. C., et al. (2014). "Single-cell trajectory detection uncovers progression and regulatory coordination in human B cell development." Cell **157**(3): 714-725.

Bendall, S. C., et al. (2011). "Single-cell mass cytometry of differential immune and drug responses across a human hematopoietic continuum." Science **332**(6030): 687-696.

Berg, D. A., et al. (2015). "Tbr2-expressing intermediate progenitor cells in the adult mouse hippocampus are unipotent neuronal precursors with limited amplification capacity under homeostasis." Frontiers in Biology **in press**.

Bodenmiller, B., et al. (2012). "Multiplexed mass cytometry profiling of cellular states perturbed by small-molecule regulators." Nat Biotechnol **30**(9): 858-867.

Bonaguidi, M. A., et al. (2011). "In vivo clonal analysis reveals self-renewing and multipotent adult neural stem cell characteristics." Cell **145**(7): 1142-1155.

Bracko, O., et al. (2012). "Gene expression profiling of neural stem cells and their neuronal progeny reveals IGF2 as a regulator of adult hippocampal neurogenesis." J Neurosci **32**(10): 3376-3387.

Buettner, F., et al. (2015). "Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells." Nat Biotechnol.

Cahan, P., et al. (2014). "CellNet: network biology applied to stem cell engineering." Cell **158**(4): 903-915.

Codega, P., et al. (2014). "Prospective identification and purification of quiescent adult neural stem cells from their in vivo niche." Neuron **82**(3): 545-559.

Encinas, J. M., et al. (2011). "Division-coupled astrocytic differentiation and age-related depletion of neural stem cells in the adult hippocampus." Cell Stem Cell **8**(5): 566-579.

Encinas, J. M., et al. (2006). "Fluoxetine targets early progenitor cells in the adult brain." Proc Natl Acad Sci U S A **103**(21): 8233-8238.

Guo, G., et al. (2013). "Mapping cellular hierarchy by single-cell analysis of the cell surface repertoire." Cell Stem Cell **13**(4): 492-505.

Hoppe, P. S., et al. (2014). "Single-cell technologies sharpen up mammalian stem cell research." Nat Cell Biol **16**(10): 919-927.

Knobloch, M., et al. (2013). "Metabolic control of adult neural stem cell activity by Fasn-dependent lipogenesis." Nature **493**(7431): 226-230.

Kriegstein, A. and A. Alvarez-Buylla (2009). "The glial nature of embryonic and adult neural stem cells." Annu Rev Neurosci **32**: 149-184.

Levsky, J. M., et al. (2002). "Single-cell gene expression profiling." Science **297**(5582): 836-840.

Li, B. and C. N. Dewey (2011). "RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome." BMC Bioinformatics **12**: 323.

Li, L. and H. Clevers (2010). "Coexistence of quiescent and active adult stem cells in mammals." Science **327**(5965): 542-545.

Lim, H. N. and A. van Oudenaarden (2007). "A multistep epigenetic switch enables the stable inheritance of DNA methylation states." Nat Genet **39**(2): 269-275.

Lu, R., et al. (2011). "Tracking single hematopoietic stem cells in vivo using high-throughput sequencing in conjunction with viral genetic barcoding." Nat Biotechnol **29**(10): 928-933.

Ming, G. L. and H. Song (2011). "Adult neurogenesis in the mammalian brain: significant answers and significant questions." Neuron **70**(4): 687-702.

Muramoto, T., et al. (2012). "Live imaging of nascent RNA dynamics reveals distinct types of transcriptional pulse regulation." Proc Natl Acad Sci U S A **109**(19): 7350-7355.

Novick, A. and M. Weiner (1957). "Enzyme Induction as an All-or-None Phenomenon." Proc Natl Acad Sci U S A **43**(7): 553-566.

Ozbudak, E. M., et al. (2004). "Multistability in the lactose utilization network of Escherichia coli." Nature **427**(6976): 737-740.

Pico, A. R., et al. (2008). "WikiPathways: pathway editing for the people." PLoS Biol **6**(7): e184.

Qiu, P., et al. (2011). "Extracting a cellular hierarchy from high-dimensional cytometry data with SPADE." Nat Biotechnol **29**(10): 886-891.

Raj, A., et al. (2006). "Stochastic mRNA synthesis in mammalian cells." PLoS Biol **4**(10): e309.

Ramskold, D., et al. (2012). "Full-length mRNA-Seq from single-cell levels of RNA and individual circulating tumor cells." Nat Biotechnol **30**(8): 777-782.

Saadatpour, A., et al. (2014). "Characterizing heterogeneity in leukemic cells using single-cell gene expression analysis." Genome Biol **15**(12): 525.

- Schofield, R. (1978). "The relationship between the spleen colony-forming cell and the haemopoietic stem cell." Blood Cells **4**(1-2): 7-25.
- Schwanhaussner, B., et al. (2011). "Global quantification of mammalian gene expression control." Nature **473**(7347): 337-342.
- Shalek, A. K., et al. (2013). "Single-cell transcriptomics reveals bimodality in expression and splicing in immune cells." Nature **498**(7453): 236-240.
- Signer, R. A., et al. (2014). "Haematopoietic stem cells require a highly regulated protein synthesis rate." Nature **509**(7498): 49-54.
- Simons, B. D. and H. Clevers (2011). "Strategies for homeostatic stem cell self-renewal in adult tissues." Cell **145**(6): 851-862.
- Stegle, O., et al. (2015). "Computational and analytical challenges in single-cell transcriptomics." Nat Rev Genet **16**(3): 133-145.
- Takeda, N., et al. (2011). "Interconversion between intestinal stem cell populations in distinct niches." Science **334**(6061): 1420-1424.
- Thorndike, R. L. (1953). "Who belongs in the family?" Psychometrika **18**(4): 267-276.
- Trapnell, C., et al. (2014). "The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells." Nat Biotechnol **32**(4): 381-386.
- Trapnell, C., et al. (2010). "Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation." Nat Biotechnol **28**(5): 511-515.
- Treutlein, B., et al. (2014). "Reconstructing lineage hierarchies of the distal lung epithelium using single-cell RNA-seq." Nature **509**(7500): 371-375.
- Usoskin, D., et al. (2015). "Unbiased classification of sensory neuron types by large-scale single-cell RNA sequencing." Nat Neurosci **18**(1): 145-153.
- Wagner, G. P., et al. (2012). "Measurement of mRNA abundance using RNA-seq data: RPKM measure is inconsistent among samples." Theory Biosci **131**(4): 281-285.

Curriculum Vitae

JAEHOON SHIN

Graduate student in the laboratory of Hongjun Song

Institute for Cell Engineering

Johns Hopkins Medical Institutions

Phone: (443) 676-5582 · E-mail: shin@jhmi.edu

Education

BS in Chemical Engineering (*summa cum laude*), Seoul National University, **2000 - 04**

MD, Seoul National University, **2004 - 08**

PhD, Johns Hopkins Medical Institutions, **2010 - current**

Scholarship, Awards, Honors

1. Honored Student, Department of Chemical Engineering, Seoul National University, Korea, **2001-04**
2. General Electric - Fulbright Scholarship, Korea, **2002-04**
3. Alumni prize, Department of Chemical Engineering, Seoul National University, Korea, **2004**
4. Doosan Scholarship, Yeonkang Foundation, Korea, **2006-08**
5. Samsung Scholarship, Samsung Foundation of Culture, Korea, **2010-2015**
6. Finalist, Innocentive challenge for single cell analysis, National Institute of Health (NIH), USA, **2015**
- NIH sponsored national competition (<http://commonfund.nih.gov/singlecell/challenge>)

Publications

1. Shin, J., Berg, D.A., Zhu, Y., Shin, J.Y., Song, J., Bonaguidi, M.A., Enikolopov, G.E., Nauen, D.W., Christian, K.M., Ming, G.-I., Song, H. (2015). Single-cell RNA-seq with Waterfall Reveals Molecular Cascades underlying Adult Neurogenesis. *Cell Stem Cell* 17(3), 360-72. (*featured as a cover*)
2. Zeng, Y., Yao, B., Shin, J., Lin, L., Kim, N., Song, Q., Liu, S., Su, Y., Guo, J., Huang, L., Wan J., Wu H., Qian J., Cheng X., Zhu H., Ming, G.-I., Jin P., Song H. (2015). Lin28A Binds Active Promoters and Recruits Tet1 to Regulate Gene Expression. *Molecular cell* 61, 1153-1160.
3. Shin, J., Ming, G.-I., Song, H. (2015). Molecular Toggle Switch of Histone Demethylase LSD1. *Molecular Cell* 57, 949-950.
4. Shin, J., Ming, G.-I., Song, H. (2015). Seeking a Roadmap toward Neuroepigenetics. *Neuron* 86, 12-15.
5. Yu, H., Su, Y., Shin, J., Zhong, C., Guo, J.U., Weng, Y.-L., Gao, F., Geschwind, D.H., Coppola, G., Ming, G.-I. (2015). Tet3 regulates synaptic transmission and homeostatic plasticity via DNA oxidation and repair. *Nature Neuroscience*.
6. Shin, J., Ming, G.-I., Song, H. (2014). Decoding neural transcriptomes and epigenomes via high-throughput sequencing. *Nature Neuroscience*, 17(11), 1463-1475.
7. Wen, Z., Nguyen, H.N., Guo, Z., Lalli, M. A., Wang, X., Su, Y., Kim, N.-S., Yoon, K.-J., Shin, J., Zhang, C., Makri, G., Nauen, D.W., Yu, H., Guzman, E., Chiang, C.H., Yoritomo, N., Kaibuchi, K., Zou, J., Christian, K.M., Cheng, L., Ross, C.A., Margolis R.L., Chen, G., Kosik, K.S., Song, H., Ming, G.-I. (2014). Synaptic dysregulation in a human iPS cell model of mental disorders. *Nature*, 515(7527), 414-418.
8. Guo, J.U., Su, Y., Shin, J.H., Shin, J., Li, H., Xie, B., Zhong, C., Hu, S., Le, T., Fan, G., Zhu, H., Chang, Q., Gao, Y., Ming, G.-I., Song, H. (2014). Distribution, recognition and regulation of non-CpG methylation in the adult mammalian brain. *Nature Neuroscience*, 17(2), 215-222.
9. Shin, J., Ming, G.-I., Song, H. (2014). DNA modifications in the mammalian brain. *Philosophical*

- Transactions of the Royal Society B: Biological Sciences, 369(1652), 20130512.
10. Hu, S., Wan, J., Su, Y., Song, Q., Zeng, Y., Nguyen, H. N., Shin, J., Cox, E., Rho, H.S., Woodard, C., Xia, S., Liu, S., Lyu, H., Ming, G.-I., Wade, H., Song, H., Qian, J., Zhu, H. (2013). DNA Methylation Presents Distinct Binding Sites for Human Transcription Factors. *eLife*, 2, e00726. doi:10.7554/eLife.00726.
 11. Weng, Y.L., An, R., Shin, J., Song, H., Ming, G.-I. (2013). DNA modifications and neurological disorders. *Neurotherapeutics*, 10(4), 556-567.
 12. Shin, J., Ming, G.-I., Song, H. (2013). By Hook or by Crook: Multifaceted DNA-Binding Properties of MeCP2. *Cell*, 152(5), 940-942.
 13. Bhang, S.H., Park, J., Yang, H.S., Shin, J., Kim, B.-S. (2013). Platelet-rich plasma enhances the dermal regeneration efficacy of human adipose-derived stromal cells administered to skin wounds. *Cell Transplantation*, 22(3), 437-445.
 14. Yoon, H.H., Bhang, S.H., Shin, J.Y., Shin, J., Kim, B.-S. (2012). Enhanced cartilage formation via three-dimensional cell engineering of human adipose-derived stem cells. *Tissue Engineering Part A*, 18(19-20), 1949-1956.
 15. *Yang, H.S., *Shin, J., Bhang, S.H., Shin, J.Y., Park, J., Im, G.I., Kim, C.S., Kim, B.-S. (2011). Enhanced skin wound healing by a sustained release of growth factors contained in platelet-rich plasma. *Experimental and Molecular Medicine*, 43(11), 622-629. (* equal contributions)

Presentations

Oral presentation

1. Shin, J. (2014). Single Cell RNA-sequencing Reveals Molecular Dynamics of Adult Hippocampal Neural Stem Cells, Oral Presentation presented at: Gordon Research Conference; Boston, MA

Poster presentation

1. Shin, J., Berg, D.A., Zhu, Y., Shin, J.Y., Nauen, D.W., Song, J., Bonaguidi, M.A., Ming, G.-I., Song, H. (2014). Single Cell RNA-seq Reveals Molecular Dynamics of Adult Hippocampal Neural Stem Cells, Poster presented at: Howard Hughes Medical Institute, Janelia Farm Conference; Ashburn, MD